



A High performance OCR System For Medieval Documents
(e.g. “Narrenschiffe” Novels of 15th Century)

Dr. Saqib Bukhari & Prof. Dr. Andreas Dengel

What we were cable of before the start of this project?

What we were cable of before the start of this project?

We have been developing  for a decade



What we were cable of before the start of this project?

We have been developing  for a decade

Google Sponsored Project to 

"A project to develop advanced OCR technologies in DFK", Google Code Blog, 2007

Complete OCR Work-Flow for Document Analysis

Breuel: *The OCROPUS Open Source OCR System*. Proceedings SPIE 20th Annual Symposium 2008



Preprocessing, Layout Analysis, OCR, ...

- Shafait, Keysers, Breuel: Performance Evaluation and Benchmarking of Six-page Segmentation Algorithms. IEEE TPAMI, 2008.
- Bukhari, Shafait, Breuel: High Performance Layout Analysis of Arabic and Urdu Document Images. IEEE ICPR, 2011.

What we were cable of before the start of this project?

We have been developing  for a decade

 Sponsored Project to 

"A project to develop advanced OCR technologies in DFK", Google Code Blog, 2007

Complete OCR Work-Flow for Document Analysis

Breuel: *The OCROPUS Open Source OCR System*. Proceedings SPIE 20th Annual Symposium 2008



Preprocessing, Layout Analysis, OCR, ...

- Shafait, Keysers, Breuel: Performance Evaluation and Benchmarking of Six-page Segmentation Algorithms. IEEE TPAMI, 2008.
- Bukhari, Shafait, Breuel: High Performance Layout Analysis of Arabic and Urdu Document Images. IEEE ICPR, 2011.

OCR Processing for Historical Documents

Breuel, Ul-Hasan, Azawi, Shafait: High Performance OCR for Printed English and Fraktur using LSTM Networks. IEEE ICDAR, 2013.

What we were cable of before the start of this project?

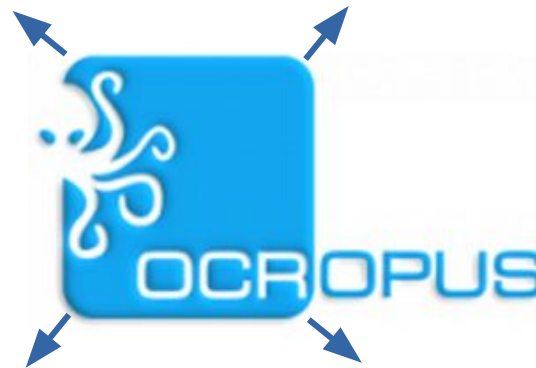
We have been developing  for a decade

 Sponsored Project to 

"A project to develop advanced OCR technologies in DFK", Google Code Blog, 2007

Complete OCR Work-Flow for Document Analysis

Breuel: *The OCROPUS Open Source OCR System*. Proceedings SPIE 20th Annual Symposium 2008



Preprocessing, Layout Analysis, OCR, ...

- Shafait, Keysers, Breuel: Performance Evaluation and Benchmarking of Six-page Segmentation Algorithms. IEEE TPAMI, 2008.
- Bukhari, Shafait, Breuel: High Performance Layout Analysis of Arabic and Urdu Document Images. IEEE ICPR, 2011.

OCR Processing for Historical Documents

Breuel, UI-Hasan, Azawi, Shafait: High Performance OCR for Printed English and Fraktur using LSTM Networks. IEEE ICDAR, 2013.

So far six PhD theses have been completed on OCROPUS, and three are on the way!

What we were cable of before the start of this project?

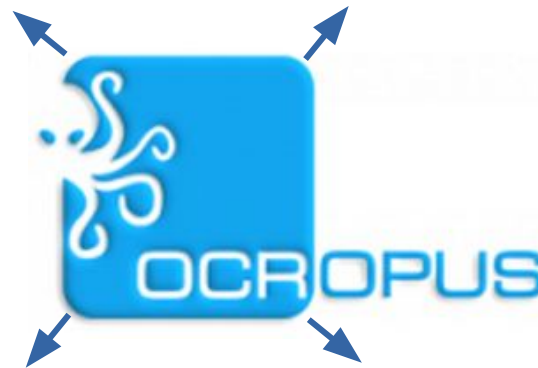
We have been developing  for a decade

 Sponsored Project to 

"A project to develop advanced OCR technologies in DFK", Google Code Blog, 2007

Complete OCR Work-Flow for Document Analysis

Breuel: *The OCROPUS Open Source OCR System*. Proceedings SPIE 20th Annual Symposium 2008



Preprocessing, Layout Analysis, OCR, ...

- Shafait, Keysers, Breuel: Performance Evaluation and Benchmarking of Six-page Segmentation Algorithms. IEEE TPAMI, 2008.
- Bukhari, Shafait, Breuel: High Performance Layout Analysis of Arabic and Urdu Document Images. IEEE ICPR, 2011.

Ocrosic

Springmann: Ocrocis, a project manager for Ocropus, Ludwig-Maximilians-University, Munich, 2015.

OCR Processing for Historical Documents

Breuel, Ul-Hasan, Azawi, Shafait: High Performance OCR for Printed English and Fraktur using LSTM Networks. IEEE ICDAR, 2013.

So far six PhD theses have been completed on OCROPUS, and three are on the way!

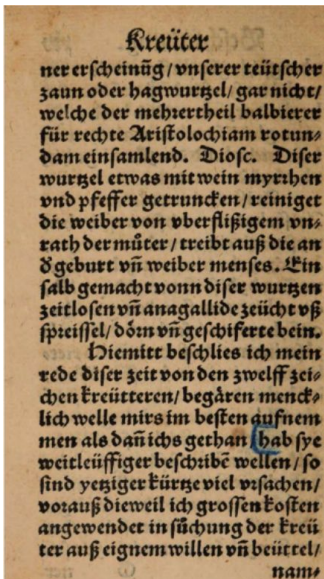
What we were cable of before the start of this project?

An Example: The



Work-Flow For 15th Century Historical Documents

Adam von Bodenstein (1557) [1]

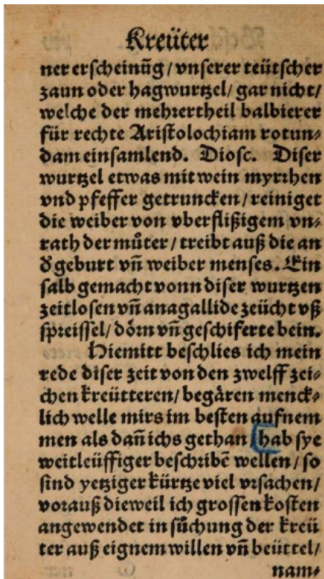


What we were cable of before the start of this project?

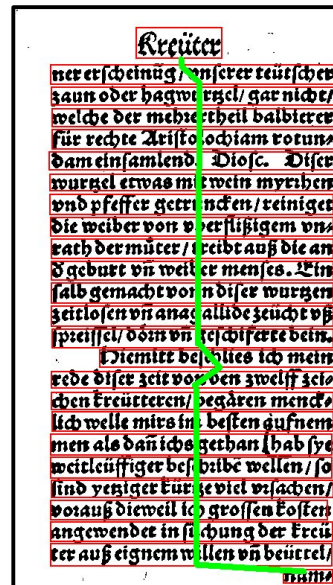
An Example: The Work-Flow For 15th Century Historical Documents



Adam von Bodenstein (1557) [1]



Preprocessing and Layout Analysis

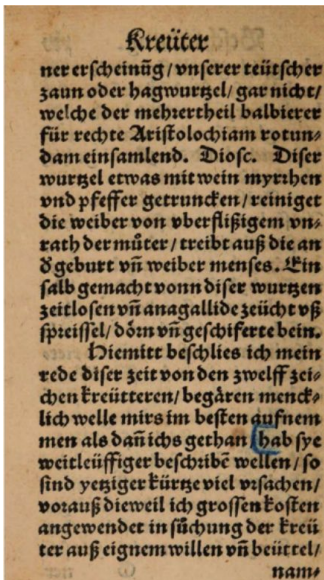


What we were cable of before the start of this project?

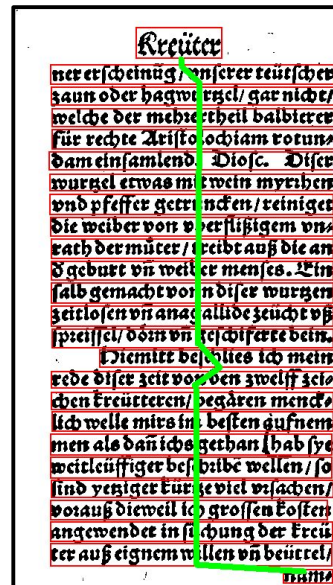
An Example: The Work-Flow For 15th Century Historical Documents



Adam von Bodenstein (1557) [1]



Preprocessing and Layout Analysis



Trained OCR Line Recognizer

Kreüter
ner erscheinüg/ vnserer teütscher
zaun oder hagwurtzel/ gar nicht/
welche der mehrertheil balbierer
für rechte Aristolochiam rotun-
dam einsamlend. Diofc. Difer
wurtzel etwas mit wein myrrhen
vnd pfeffer getruncken/ reiniget
die weiber von vberflüzigem vn-
rath der müter/ treibt auß die an-
d geburt vñ weiber menses. Ein
salb gemacht vonn difer wurzen
zeitlosen vñ anagallide zeücht vñ
spreißel/ dönn vñ geschiferte bein.
Hiemitt beschlies ich mein
rede difer zeit von den zwelff ze-
chen kreütteren/ begären menck-
lich welle mirs im besten aufnem-
men als daß ichs gethan/ hab sye
weitleüffiger beschribē wellen/ so
sind yetziger kürzte viel vrsachen/
voraus dieweil ich groffen kosten
angewendet in süchung der kreü-
ter auß eignem willen vñ beütel/
nen:

OCROPUS Results [1]

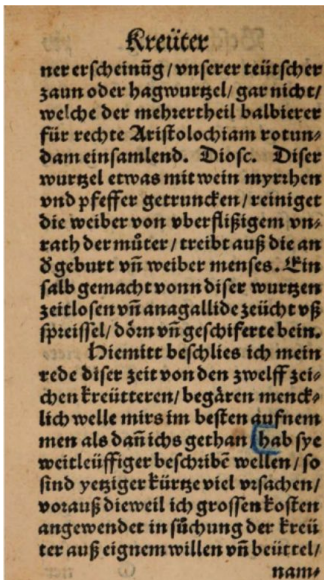
What we were cable of before the start of this project?

An Example: The Work-Flow For 15th Century Historical Documents

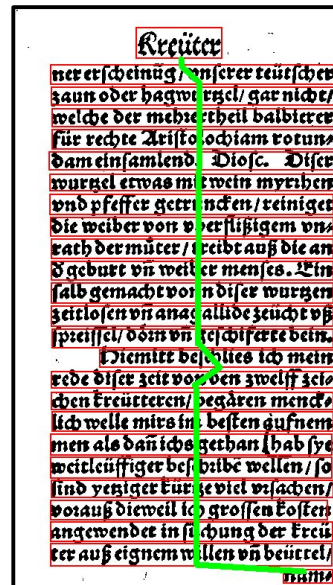


On Test images from Adam von Bodenstein (1557), the performance:
Ocropus 99%, **ABBY 85%** and **Tesseract 78%** [1].

Adam von Bodenstein (1557) [1]



Preprocessing and Layout Analysis



Trained OCR Line Recognizer

Kreüter
ner erscheinüg/ vnserer teütscher
zaun oder hagwurtzel/ gar nicht/
welche der mehrertheil baldieret
für rechte Aristolochiam rotun-
dam einsamlend. Diosc. Difer
wurtzel etwas mit wein myrrhen
vnd pfeffer getruncken/ reiniget
die weiber von vberflüzigem vn-
rath der müter/ treibt auß die an-
geburt vn weiber menses. Ein
salb gemacht vonn diser wurtzen
zeitlosen vn anagallide zeucht vß
spreißel/ dönn vn geschiferte bein.
Hiemitt beschlies ich mein
rede diser zeit von den zwelff ze-
chen kreütteren/ begären menck-
lich welle mirs im besten aufnem-
men als daß ichs gethan/ hab fye
weitleuffiger beschribē wellen/ so
sind yetziger kürze viel vrsachen/
voraus dieweil ich groffen kosten
angewendet in führung der kreü-
ter auß eignem willen vn beütel/
nen:

OCROPUS Results [1]

Kreücev
ner erscheinung / vnserer teutscher
zäun oderhagwuryel/ garnicht/
welche der mehrertheil baldieret
für rechte Aristolochiam rotun-
dameinsamlend. Diosc. Diser
würget etwas mit wein Myrrhen
vnd Pfeffer getrvncken/ reiniget
die weider von vderssissigem vn-
rath der müter/ treibt auß die an-
Lgedurt vn weider menses. Ein
faldgemachtvonndiser würgen
zeitlosenvnanagallidezeüchtvß
spreissel/ dönnvngeschiferte dein.
iZiemitt deschlies ich mein
rede diser zeitvonden zwelffzei-
chenkreütteren/degären menck-
ltch welle mirs im dessen aufnem-
men als dan ichs gethan / had fye
weitleufftgerdeschride wellen / so
sind yegigerkürgeviel vrsachen/
vorausdieweil ichgrossinkosten
angewendet insüchungder kreü-
ter auß eignem willen vndeüerel/
nen:

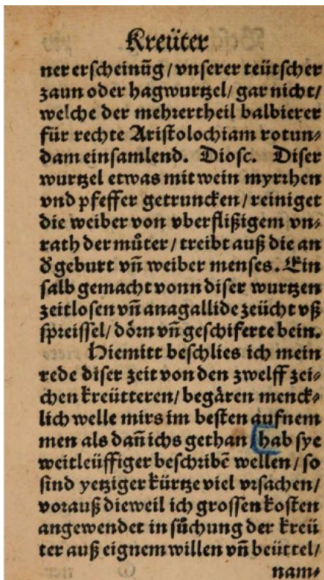
ABBY Results [1]

What we were cable of before the start of this project?

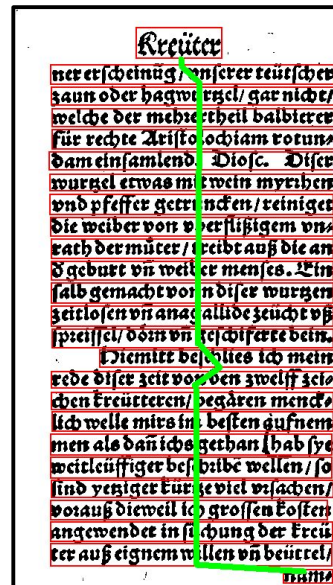
An Example: The Work-Flow For 15th Century Historical Documents

On Test images from Adam von Bodenstein (1557), the performance:
Ocropus 99%, **ABBY 85%** and **Tesseract 78%** [1].

Adam von Bodenstein (1557) [1]



Preprocessing and Layout Analysis



Trained OCR Line Recognizer

Kreüter
ner erscheinüg/ vnserer teütscher
zaun oder hagwurtzel/ gar nicht/
welche der mehrertheil baldieret
für rechte Aristolochiam rotun-
dam einsamlend. Diosc. Difer
wurtzel etwas mit wein myrrhen
vnd pfeffer getruncken/ reiniget
die weiber von vberflüzigem vn-
rath der müter/ treibt auß die an-
geburt vn weiber mensces. Ein
salb gemacht vonn diser wurtzen
zeitlosen vn anagallide zeucht vß
spreißel/ dönn vn geschiferte bein.
Hiemitt beschlies ich mein
rede diser zeit von den zwelff ze-
chen kreütteren/ begären menck-
lich welle mirs im besten aufnem-
men als daß ichs gethan/ hab fye
weitleuffiger beschribē wellen/ so
sind yetziger kürze viel vrsachen/
voraus dieweil ich groffen kosten
angewendet in führung der kreü-
ter auß eignem willen vn beütet/
nen

Ocropus Results [1]

Kreücev
ner erscheinung / vnserer teutscher
zäun oderhagwuryel/ garnicht/
welche der mehrertheil baldieret
für rechte Aristolochiam rotun-
dameinsamlend. Diosc. Diser
würget etwas mit wein Myrrhen
vnd Pfeffer getrvncken/ reiniget
die weiber von vderssissigem vn
rath dermäter/ treibt auß die an
Lgedurt vn weiber mensces. Ein
faldgemachtvonndiser würgen
zeitlosenvnanagallidezeüchtvß
spreissel/ dönnvngeschiferte dein.
iZiemitt deschlies ich mein
rede diser zeitvonden zwelffzei-
chenkreütteren/degären menck-
ltch welle mirs im dessen aufnem
men als dan ichs gethan / had fye
weitleüfftiger beschribē wellen/ so
sind yegigerkürzeviel vrsachen/
vorausdieweil ichgrossinkosten
angewendet insüchungder kreü
ter auß eignem willen vndeüercl/
nen

ABBY Results [1]

so what else is required for the 15th Century “Narrenschiffe” novels of this project?

What are the limitations of the OCRopus?

What are the limitations of the OCRopus?

Layout Analysis

What are the limitations of the OCRopus?

Layout Analysis

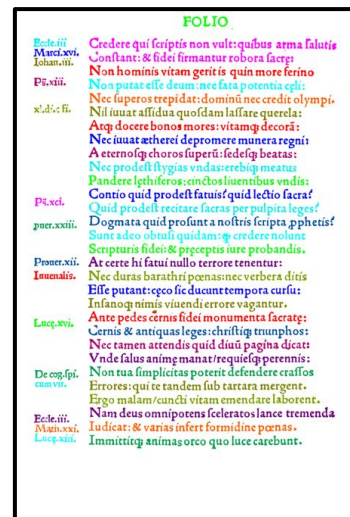
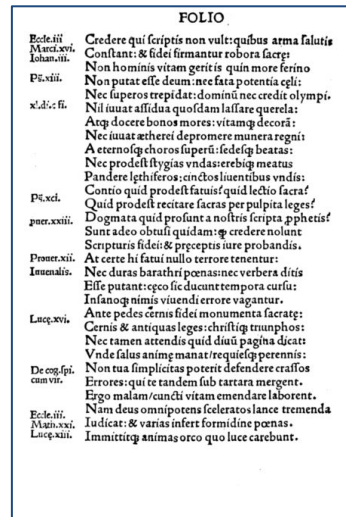
FOLIO	
Ecle.iii. Mand.xvi. Iohann.xiii.	Credere qui scriptis non vult: quibus arma salutis Constant: & fidei firmantur robora sacre:
Pet.xiii.	Non hominis vitam gerit is quin more ferino
x'.d'.c. fi.	Non putat esse deum: nec fata potentia egli: Nec superos trepidat: dominū nec credit olympi: Nil iuuat affidua quosdam lassare querela: Atq; docere bonos mores: vitamq; decorā: Nec iuuat aetherei depromere munera regni: A eternisq; choros superū: fedelq; beatas: Nec prodest itygias vndas: erebq; meatas Pandere lghiferos: cinctos luentibus vndis: Contio quid prodest fatuis: quid lectio sacra?
Pet.xci.	Quid prodest recitare sacras per pulpita leges?
grec.xxiii.	Dogmata quid profant a nostris scripta pphetis?
	Sunt adeo obtusi quidam: q; credere nolunt Scripturis fidei: & pceptis iure probandis.
Præter.xii.	At certe hi fatui nullo terrore tenentur:
Isaenalis.	Nec duras barathri poenas: nec verbera ditis Esse putant: ego sic ducunt tempora curfu: Infanoq; nimis viuendi errore vagantur.
Luce.xvi.	Ante pedes cernis fidei monumenta sacratq; Cernis & antiquas leges: christiq; triumphos: Nec tamen attendis quid diuū pagina dicat: Vnde salus animæ manat: requiesq; perennis:
De cog.fpi. can.vii.	Non tua simplicitas poterit defendere crassos Errores: qui te tandem sub tartara mergent. Ergo malam/cuncti vitam emendare laborent. Nam deus omnipotens: sceleratos lance tremenda Iudicat: & varias inferi formidine penas.
Ecle.iii. Mand.xvi. Luce.xiii.	Immittitq; animas orco quo luce carebunt.

FOLIO	
Ecle.iii. Mand.xvi. Iohann.xiii.	Credere qui scriptis non vult: quibus arma salutis Constant: & fidei firmantur robora sacre: Non hominis vitam gerit is quin more ferino
Pet.xiii.	Non putat esse deum: nec fata potentia egli: Nec superos trepidat: dominū nec credit olympi: Nil iuuat affidua quosdam lassare querela: Atq; docere bonos mores: vitamq; decorā: Nec iuuat aetherei depromere munera regni: A eternisq; choros superū: fedelq; beatas: Nec prodest itygias vndas: erebq; meatas Pandere lghiferos: cinctos luentibus vndis: Contio quid prodest fatuis: quid lectio sacra?
Pet.xci.	Quid prodest recitare sacras per pulpita leges?
grec.xxiii.	Dogmata quid profant a nostris scripta pphetis?
	Sunt adeo obtusi quidam: q; credere nolunt Scripturis fidei: & pceptis iure probandis.
Præter.xii.	At certe hi fatui nullo terrore tenentur:
Isaenalis.	Nec duras barathri poenas: nec verbera ditis Esse putant: ego sic ducunt tempora curfu: Infanoq; nimis viuendi errore vagantur.
Luce.xvi.	Ante pedes cernis fidei monumenta sacratq; Cernis & antiquas leges: christiq; triumphos: Nec tamen attendis quid diuū pagina dicat: Vnde salus animæ manat: requiesq; perennis:
De cog.fpi. can.vii.	Non tua simplicitas poterit defendere crassos Errores: qui te tandem sub tartara mergent. Ergo malam/cuncti vitam emendare laborent. Nam deus omnipotens: sceleratos lance tremenda Iudicat: & varias inferi formidine penas.
Ecle.iii. Mand.xvi. Luce.xiii.	Immittitq; animas orco quo luce carebunt.

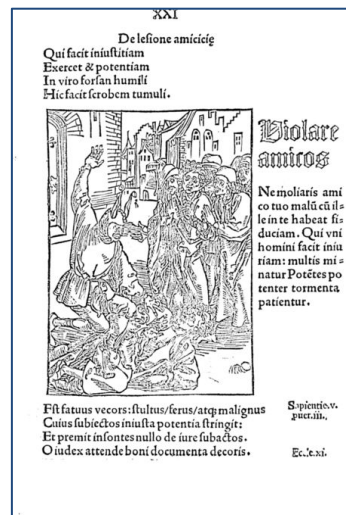
OCRopus 100%

What are the limitations of the OCRopus?

Layout Analysis



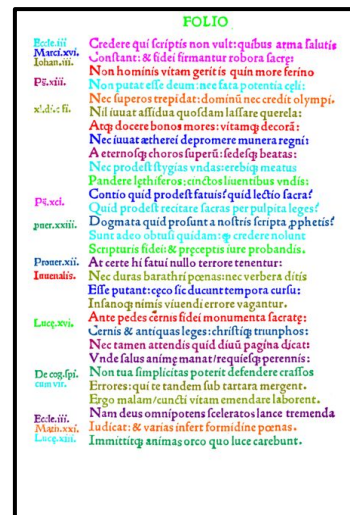
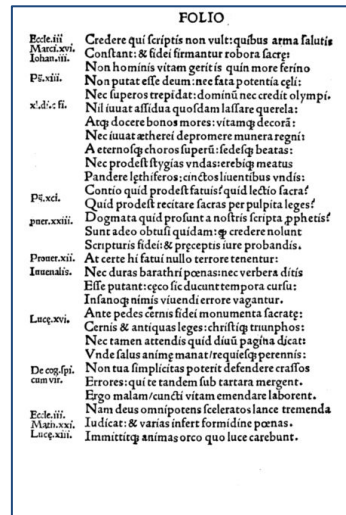
OCRopus 100%



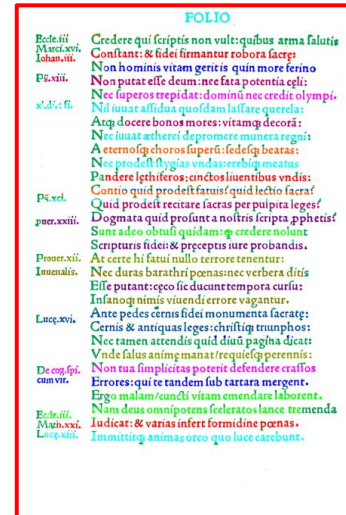
OCRopus 75%

What are the limitations of the OCRopus?

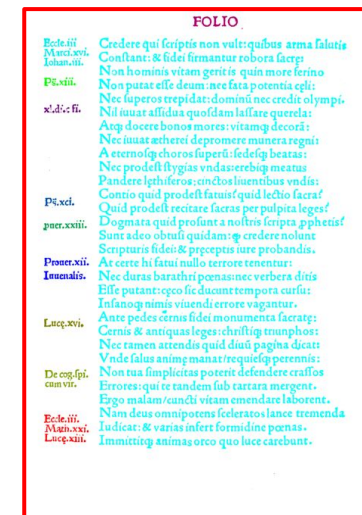
Layout Analysis



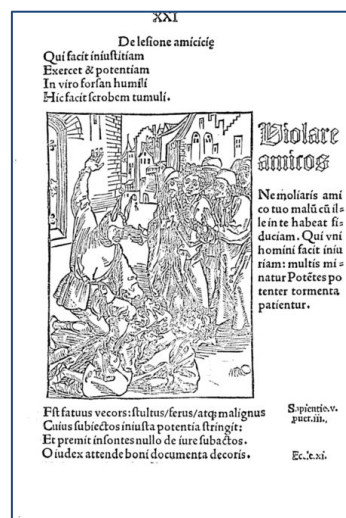
OCRopus 100%



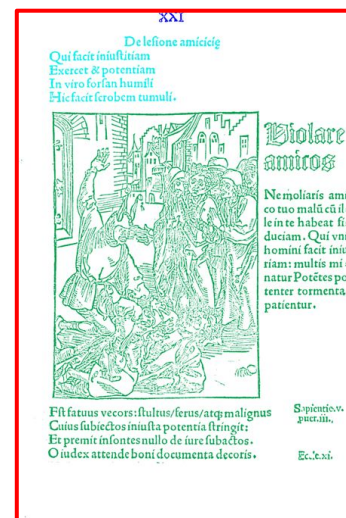
XY Cut < 5%



Voronoi < 5%



OCRopus 75%



XY Cut < 5%

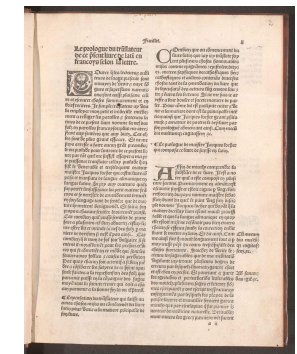
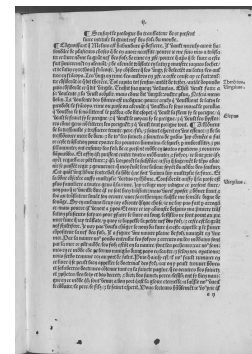
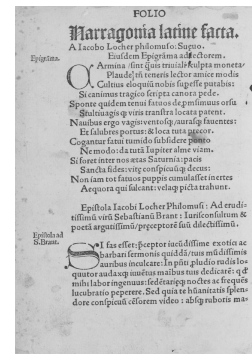
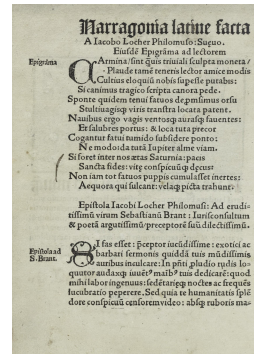
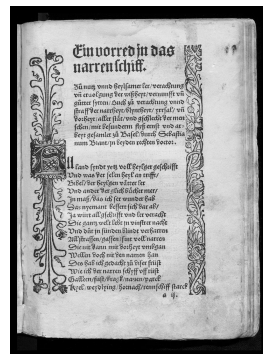
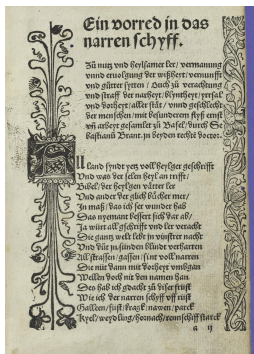
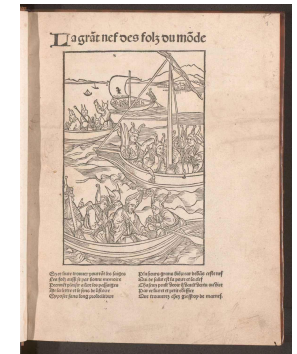
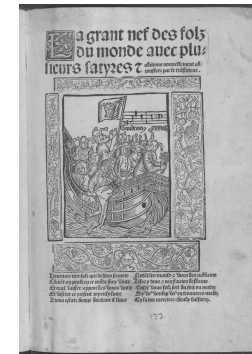
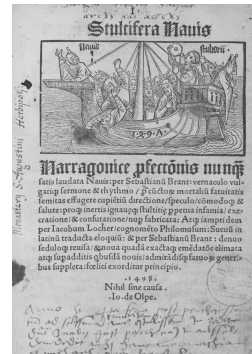
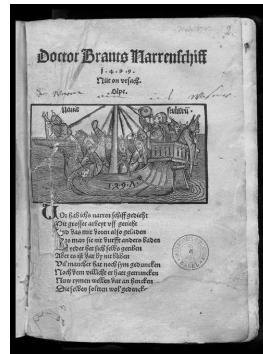


Voronoi < 5%

What are the limitations of the OCRopus?

Layout Analysis

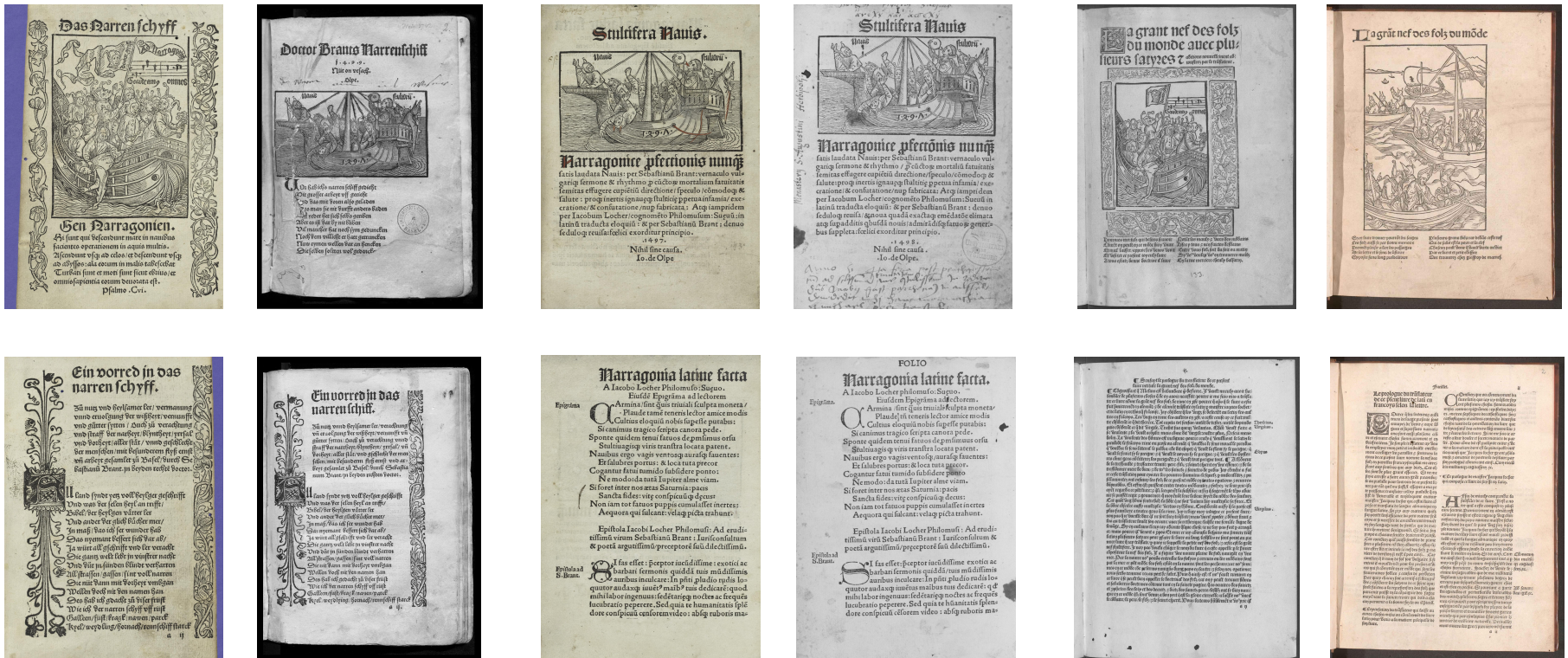
What are we dealing in this project?



What are the limitations of the OCRopus?

Layout Analysis

What are we dealing in this project?



so what else is required for the 15th Century “Narrenschiffe” novels of this project?

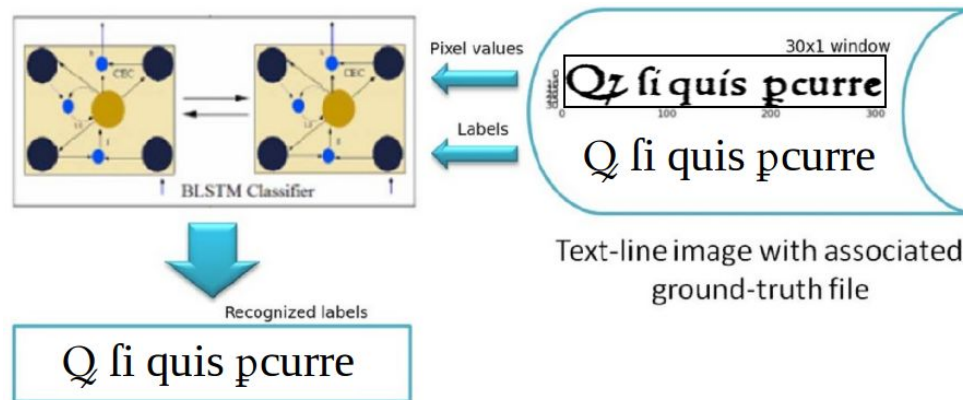
(i) Advanced Layout Analysis

What are the limitations of the OCRopus?

OCR Model

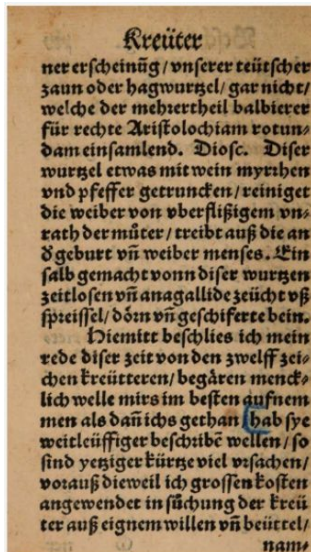
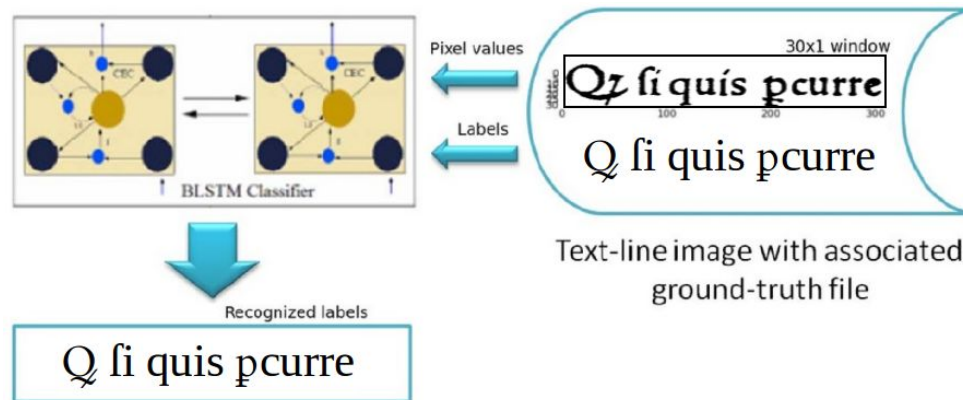
What are the limitations of the OCRopus?

OCR Model

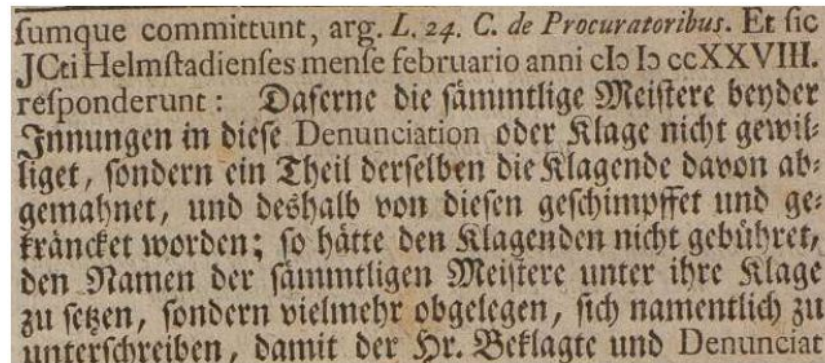


What are the limitations of the OCRopus?

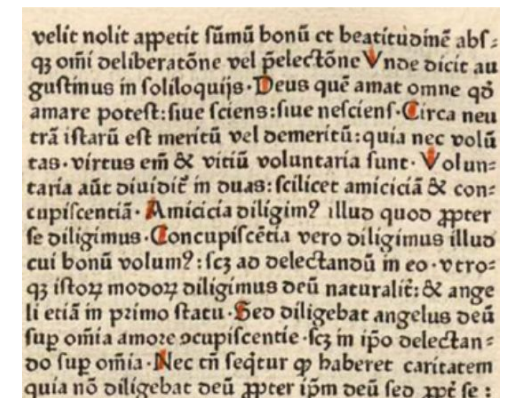
OCR Model



Adam von Bodenstein 1557 [1]
(OCRopus: 99%, ABBYY: 85%. Tesseract: 78%)



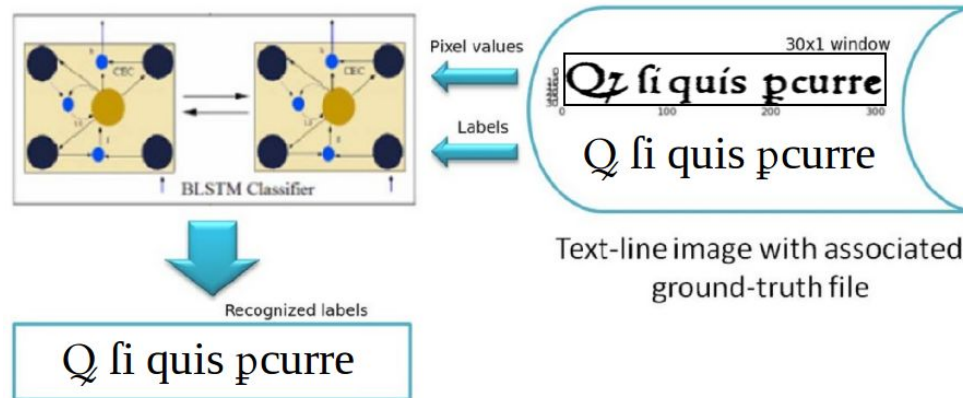
Augustinus Leyer 1735 [1]
(OCRopus: 97%, ABBYY: 77%. Tesseract: 82%)



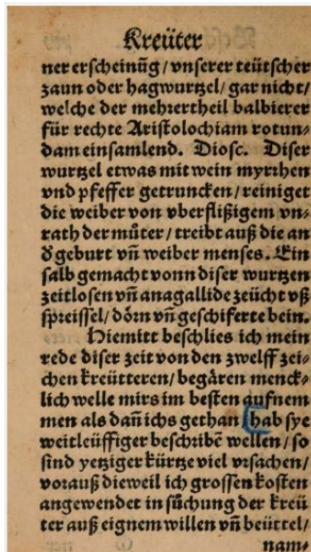
Augsburg before 1476 [1]
(OCRopus: 98%)

What are the limitations of the OCRopus?

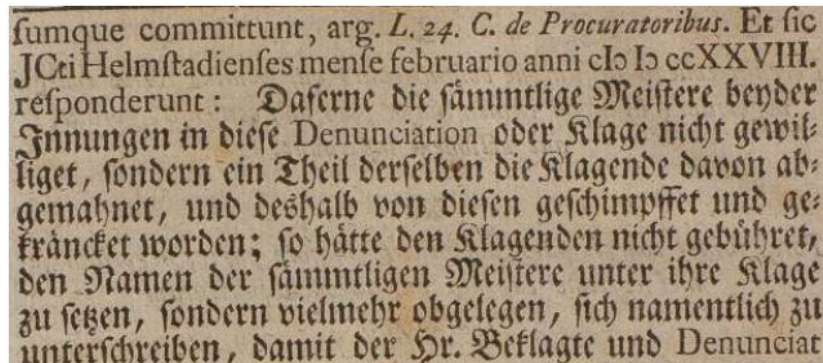
OCR Model



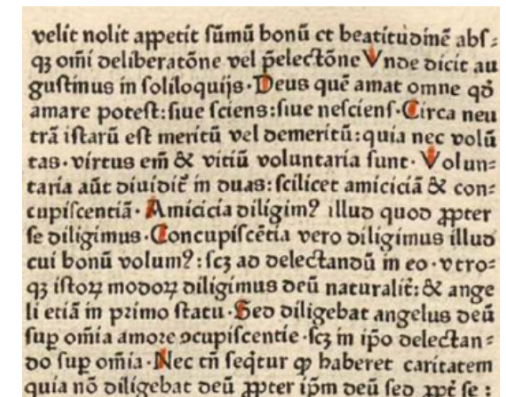
Training Data
50K Lines



Adam von Bodenstein 1557 [1]
(OCRopus: 99%, ABBYY: 85%, Tesseract: 78%)



Augustinus Leyer 1735 [1]
(OCRopus: 97%, ABBYY: 77%, Tesseract: 82%)



Augsburg before 1476 [1]
(OCRopus: 98%)

What are the limitations of the OCRopus?

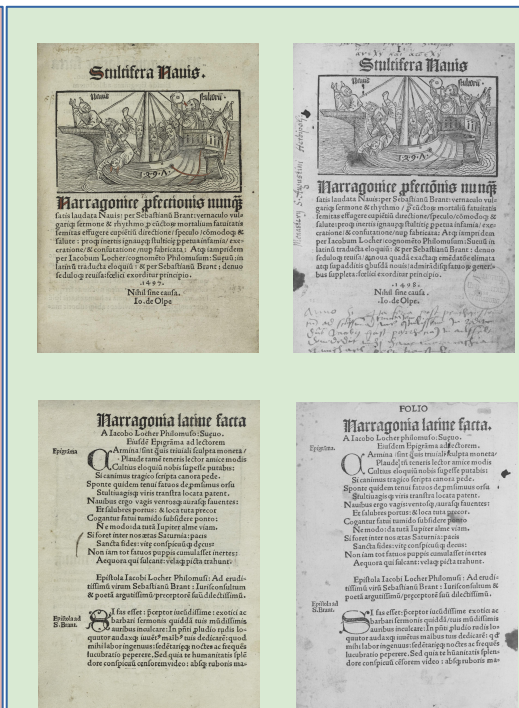
OCR Model

What are we dealing in this project?

German



Latin



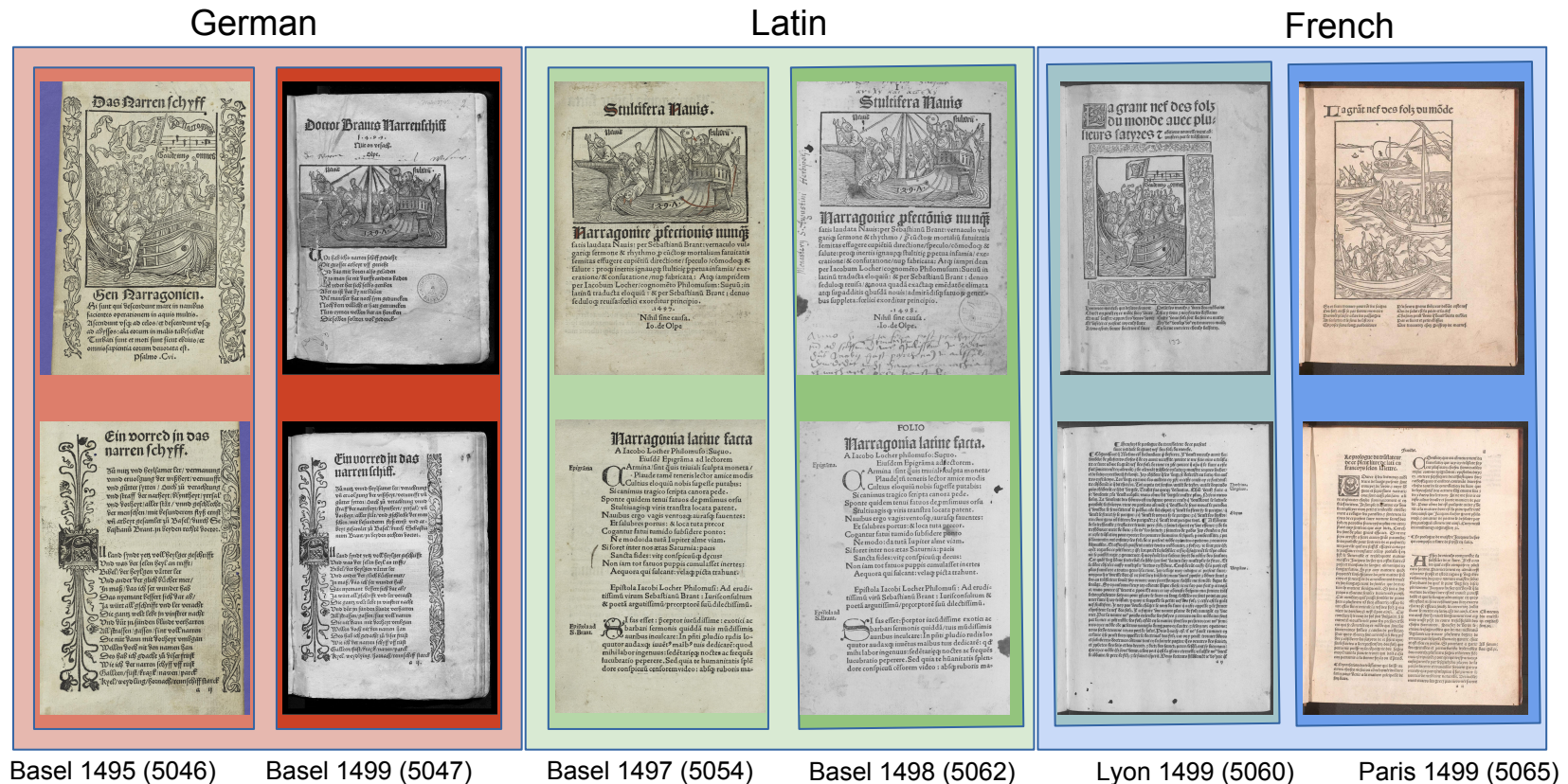
French



What are the limitations of the OCRopus?

OCR Model

What are we dealing in this project?



Basel 1495 (5046)

Basel 1499 (5047)

Basel 1497 (5054)

Basel 1498 (5062)

Lyon 1499 (5060)

Paris 1499 (5065)

What are the limitations of the OCRopus?

OCR Model

What are we dealing in this project?



so what else is required for the 15th Century “Narrenschiffe” novels of this project?

(ii) Automatic OCR Model (anyOCR!)

Our Main Goals in This Project!



Our Main Goals in This Project!

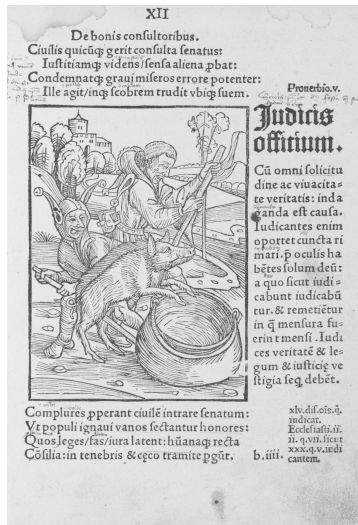
Advanced Layout Analysis

Automatic OCR Model (anyOCR)



Our Main Goals in This Project!

Advanced Layout Analysis



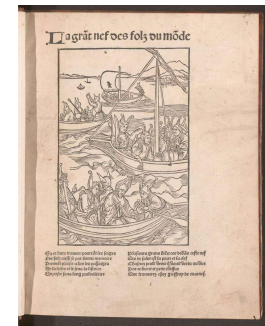
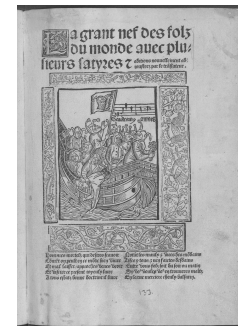
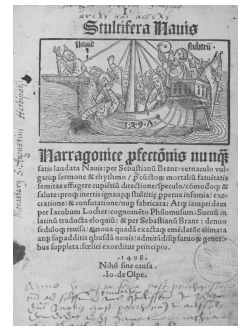
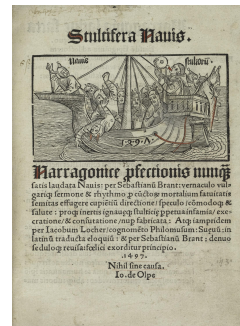
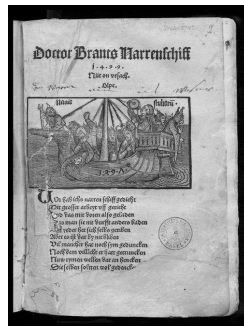
XII
De bonis confultoribus.
Ciuiilis quicūq; gerit confulta fenatus:
Iu itiamq; videns/fenfa aliena pbat:
Condemnatq; graui miferos errore poterit:
Prouerbio. v.
Ille agit/inq; scrobrem trudit vbiq; fuem.

Iudicis
offitium.
Cū omni sollicitudine ac viuacitate veritatis: inda ganda e causa. Iudicantes enim oportet cun a rimari. p oculis habētes solum deū: a quo sicut iudicabunt iudicabuntur. & remetiētur in q mensura fuerint mensi. Iudices veritatē & legum & iu icię ve ligia seq debet.

Complures pperant ciuile intrare fenatum:
xliv. dif. ois. q.
iudicat.
Ecclesia i. ii.
ii. q. vii. sicut
xxx. q. v. iudi
cantem.

Vt populi ignaui vanos se antur honores:
Quos leges/fas/iura latent: hūanaq; re a
Cōsilia: in tenebris & ceco tramite pgiūt.

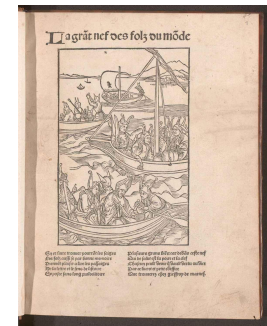
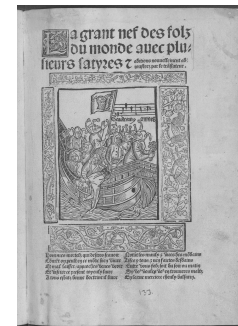
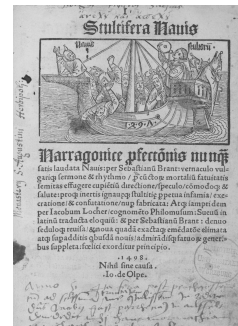
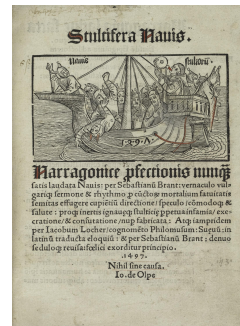
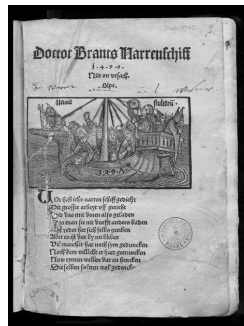
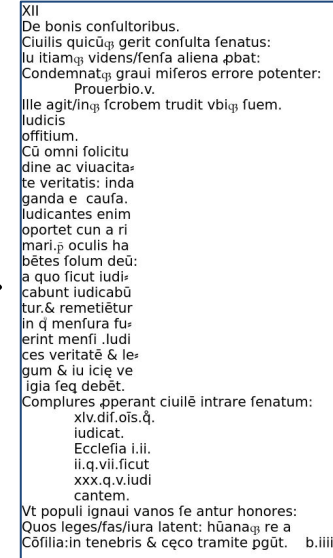
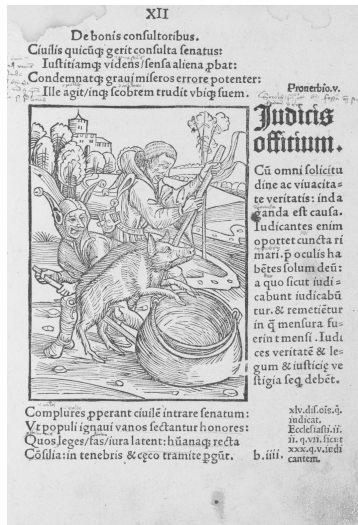
b. iiii.



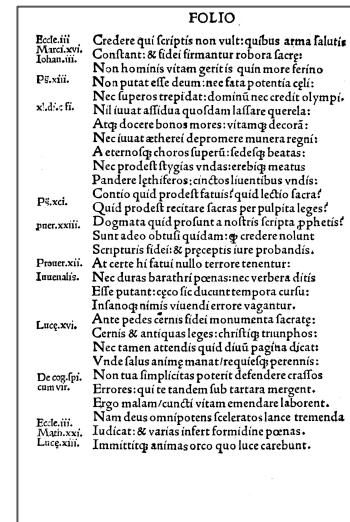
Our Main Goals in This Project!

Advanced Layout Analysis

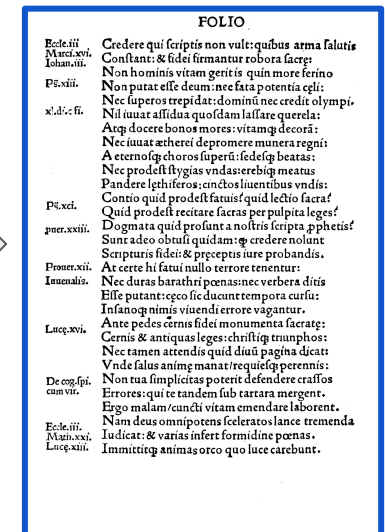
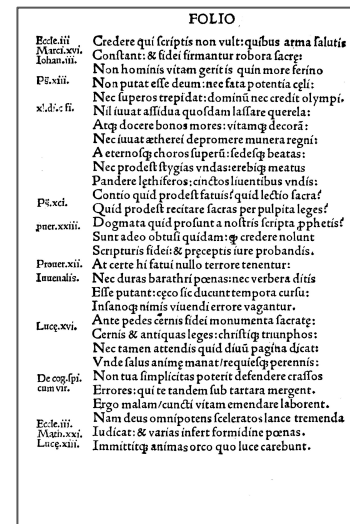
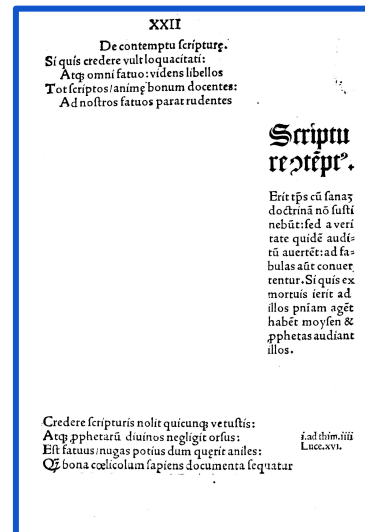
Automatic OCR Model (anyOCR)



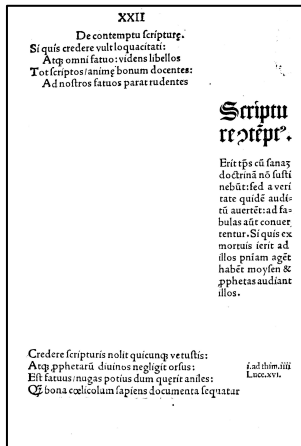
Text and Non-Text Segmentation Method



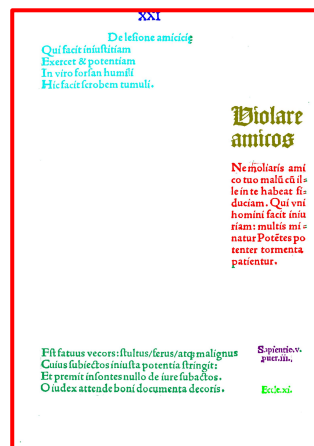
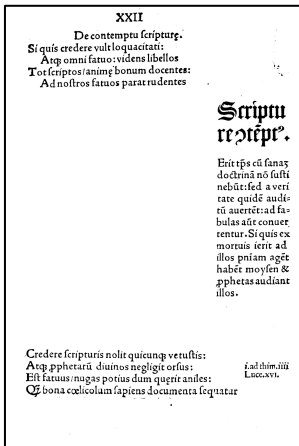
Text and Non-Text Segmentation Method



Text Line Segmentation



Text Line Segmentation



XY Cut

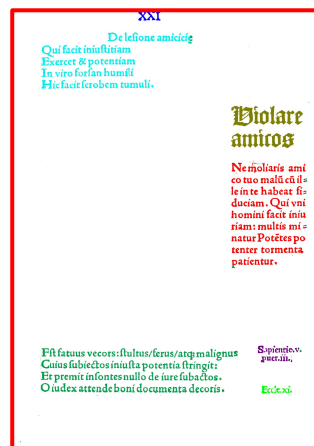
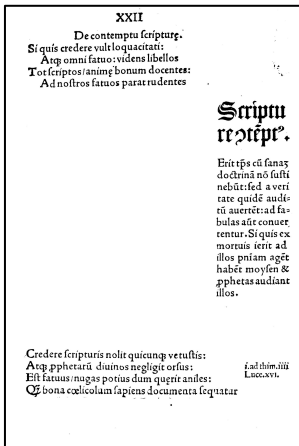


Voronoi



OCRopus

Text Line Segmentation



XY Cut



Voronoi



OCRopus



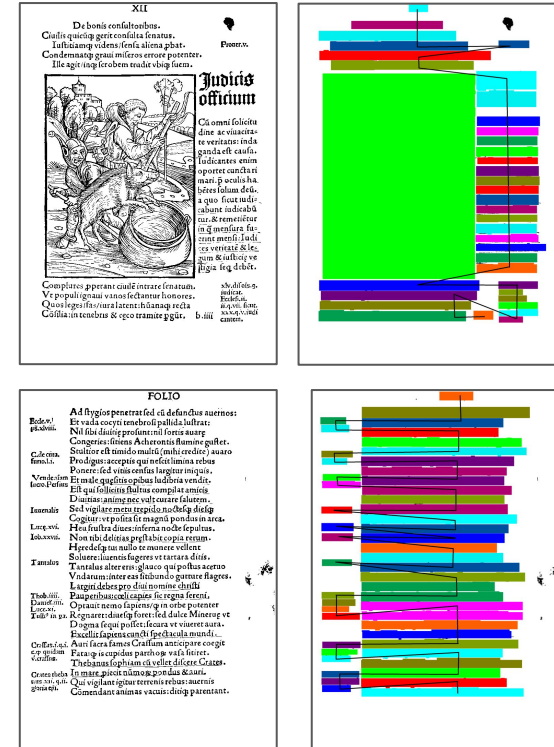
OCRopus++



Advanced Layout Analysis

Performance Evaluation of **OCROPUS++**:

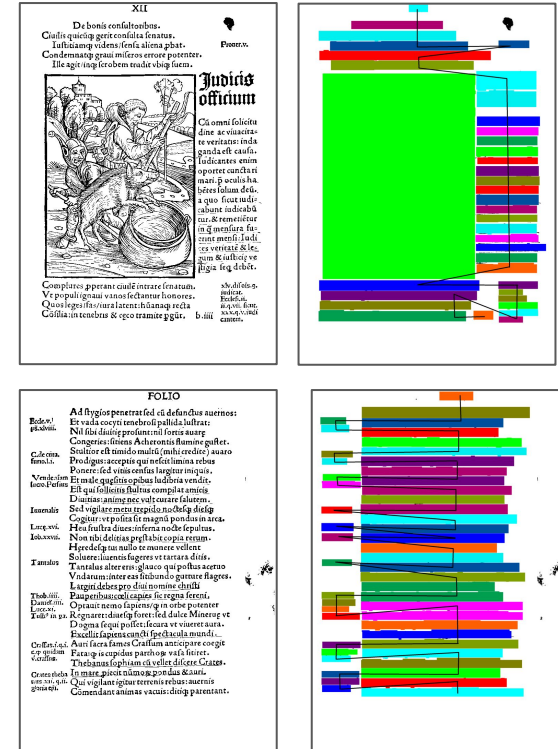
Performance Evaluation of OCROPUS++:



Performance Evaluation of **OCRopus++**:

- Text and Non-Text Segmentation Accuracy: **99.34%**
- Text Line Segmentation Accuracy: **87%***

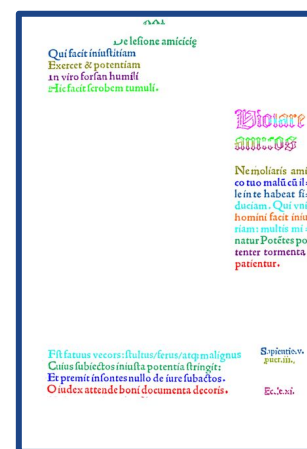
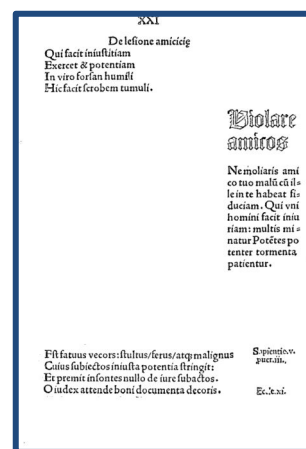
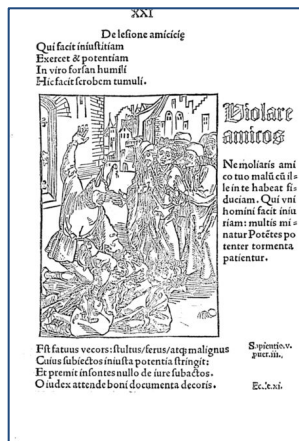
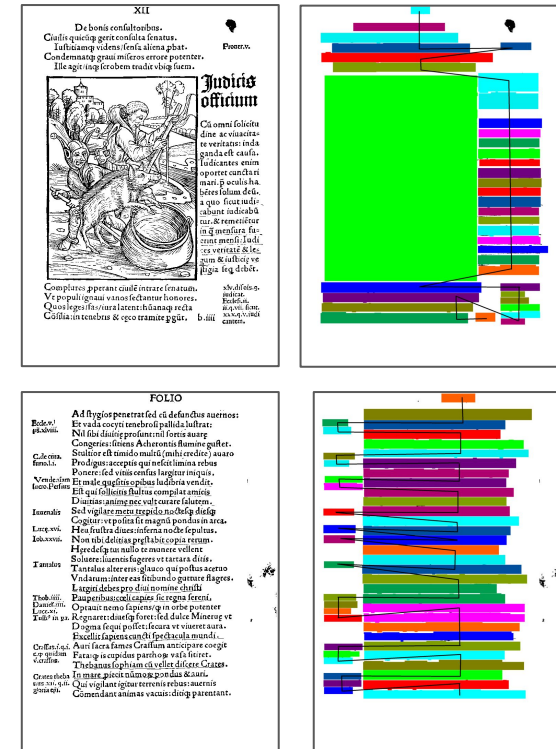
*as compared to state-of-the-art methods ~10% and **OCRopus 80%.**



Performance Evaluation of OCROPUS++:

- Text and Non-Text Segmentation Accuracy: **99.34%**
- Text Line Segmentation Accuracy: **87%***

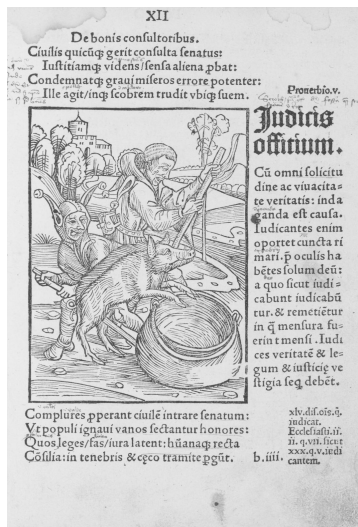
*as compared to state-of-the-art methods ~10% and OCROPUS 80%.



Our Main Goals in This Project!

Advanced Layout Analysis

Automatic OCR Model (anyOCR)

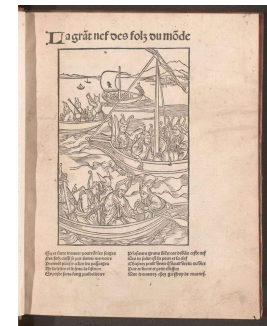
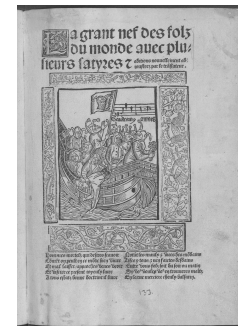
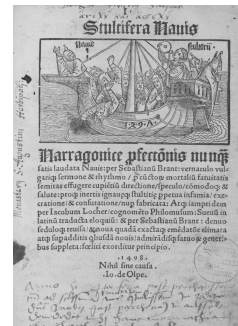
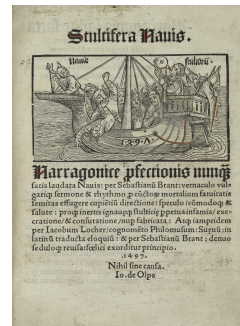
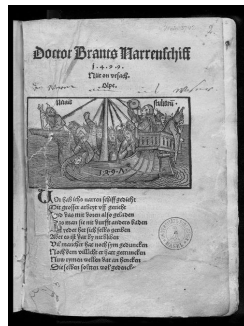


XII
De bonis confultoribus.
Ciuilis quicūq; gerit confulta fenatus:
Iu itiamq; videns/fenfa aliena pbat:
Condemnatq; graui miferos errore potenter:
Ille agit/inq; scrobrem trudit vbiq; fuem.

Judicis officium.
Cū omni sollicitudine ac viuacitate veritatis: inda ganda e causa. Iudicantes enim oportet cun a rimari. p oculis habētes solum deū: a quo sicut iudicabunt iudicabūtur. & remetiētur in q mensura fuerint mensi. Iudices veritatē & legum & iu icię ve ligia seq debēt.

Complures pperant ciuilē intrare fenatum:
xliv. dif. ois. q. iudicat.
Ecclesia i. ii.
ii. q. vii. sicut
xxx. q. v. iudi cantem.

Vt populi ignaui vanos se antur honores:
Quos leges/fas/iura latent: hūanaq; re a
Cōsilia: in tenebris & cęco tramite pğūt.





anyOCR - Automatic OCR Model

Background: OCR Training Models can broadly be classified as:

- **Segmentation-based OCR (Tesseract)**
 - individual characters classification
- **Segmentation-Free OCR (OCRopus)**
 - line recognizer

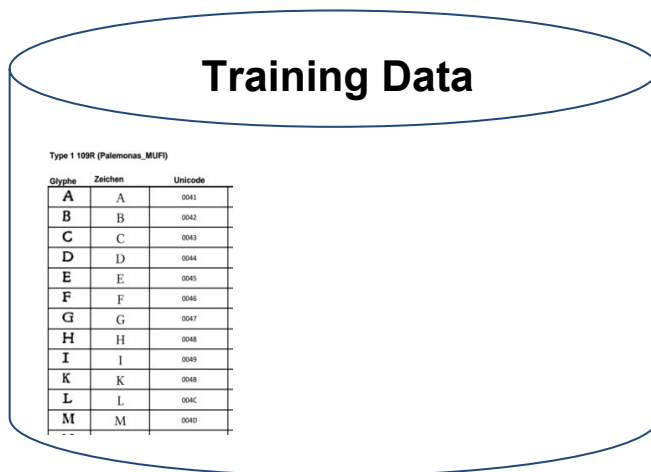


anyOCR - Automatic OCR Model

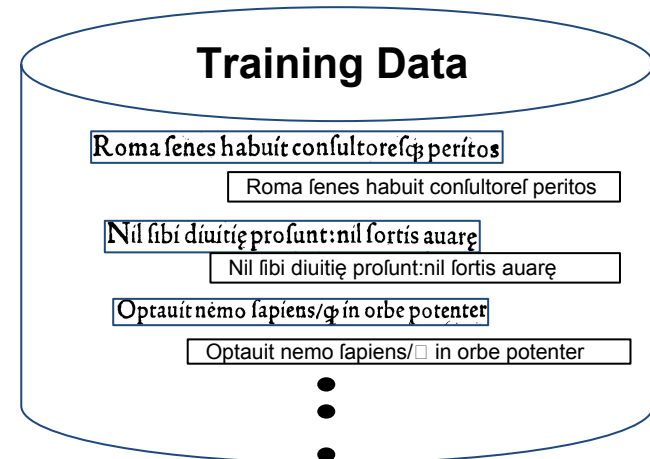
Background: OCR Training Models can broadly be classified as:

- **Segmentation-based OCR (Tesseract)**
 - individual characters classification
- **Segmentation-Free OCR (OCROPUS)**
 - line recognizer

TypeTable



Tools: Aletheia and Franken+



~50k to 100k

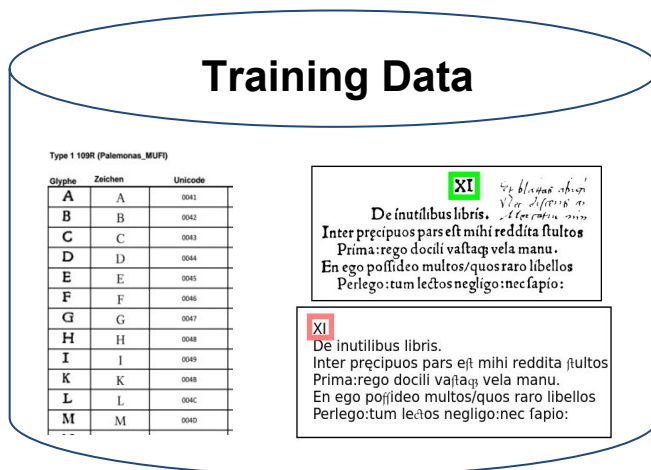


anyOCR - Automatic OCR Model

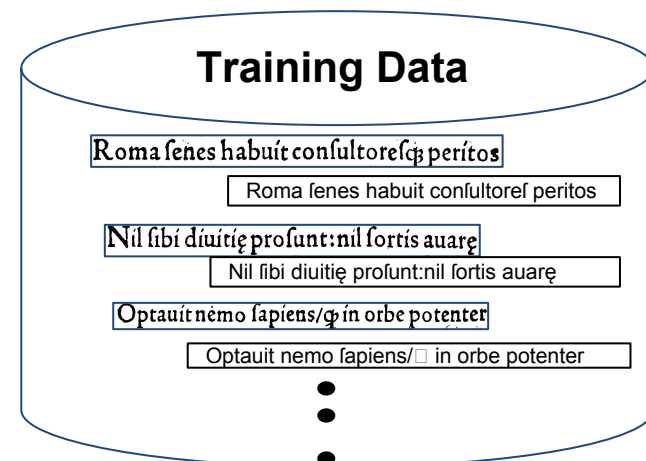
Background: OCR Training Models can broadly be classified as:

- **Segmentation-based OCR (Tesseract)**
 - individual characters classification
- **Segmentation-Free OCR (OCRopus)**
 - line recognizer

TypeTable / Few Pages



Tools: Aletheia and Franken+



~50k to 100k

Background: OCR Training Models can broadly be classified as:

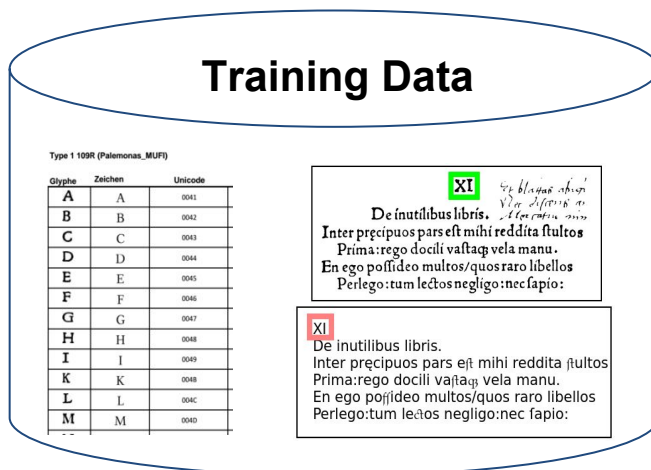
- **Segmentation-based OCR (Tesseract)**

- individual characters classification
- **Lower Performance**
 - Sensitive to Noise
 - Character Segmentation Errors
 - No Language Model

- **Segmentation-Free OCR (OCRopus)**

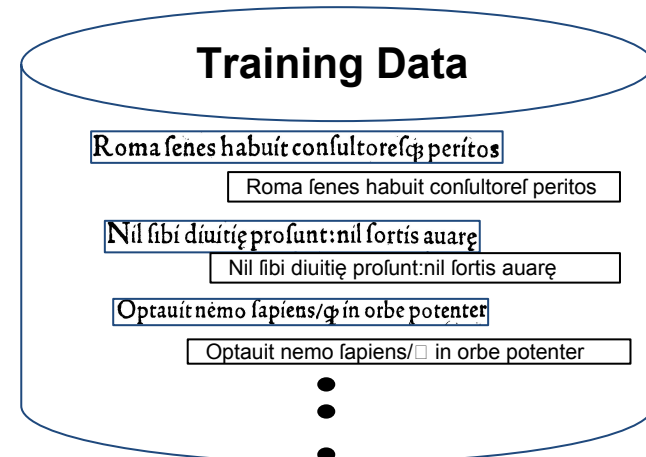
- line recognizer
- **Higher Performance**
 - Insensitive to Noise
 - Independent to Noise & Character Seg.
 - Implicit Language Model

TypeTable / Few Pages



Tools: Aletheia and Franken+

Training Data

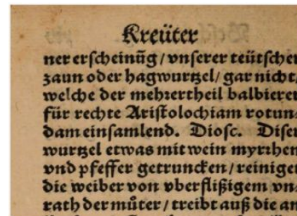


~50k to 100k

Background: OCR Training Models can broadly be classified as:

- **Segmentation-based OCR (Tesseract)**

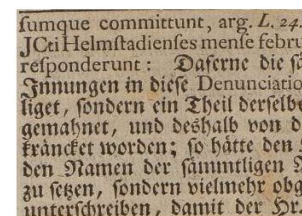
- individual characters classification
- **Lower Performance**
 - Sensitive to Noise
 - Character Segmentation Errors
 - No Language Model



Adam von Bodenstein 1557
(Ocropus: 99%, Tesseract: 78%)

- **Segmentation-Free OCR (OCRopus)**

- line recognizer
- **Higher Performance**
 - Insensitive to Noise
 - Independent to Noise & Character Seg.
 - Implicit Language Model



Augustinus Leyer 1735

(OCRopus: 97%, Tesseract: 82%)

TypeTable / Few Pages

Training Data

Type 1 109R (Palenomas_MUP)

Glyphe	Zeichen	Unicode
A	A	0041
B	B	0042
C	C	0043
D	D	0044
E	E	0045
F	F	0046
G	G	0047
H	H	0048
I	I	0049
K	K	004B
L	L	004C
M	M	004D

XI De inutilibus libris.
Inter precipuos pars est mihi reddita stultos
Prima:rego docili vastaq; vela manu.
En ego possideo multos/quos raro libellos
Perlego:tum lectos negligo:nec fapio:

XI De inutilibus libris.
Inter precipuos pars est mihi reddita stultos
Prima:rego docili vastaq; vela manu.
En ego possideo multos/quos raro libellos
Perlego:tum lectos negligo:nec fapio:

Tools: Aletheia and Franken+

Training Data

Roma fenes habuit consultoresq; peritos

Roma fenes habuit consultores peritos

Nil sibi diuitiæ profunt:nil fortis auaræ

Nil sibi diuitiæ profunt:nil fortis auaræ

Optauit nemo sapiens/ in orbe potenter

Optauit nemo sapiens/ in orbe potenter

~50k to 100k

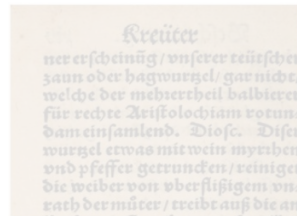
Background: OCR Training Models can broadly be classified as:

- **Segmentation-based OCR (Tesseract)**

- individual characters classification

○ Lower Performance

- Sensitive to Noise
- Character Segmentation Errors
- No Language Model



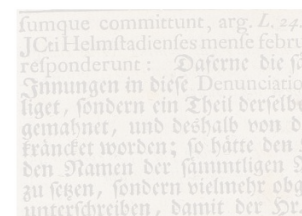
Adam von Bodenstein 1557
(Ocropus: 99%, Tesseract: 78%)

- **Segmentation-Free OCR (OCRopus)**

- line recognizer

○ Higher Performance

- Insensitive to Noise
- Independent to Noise & Character Seg.
- Implicit Language Model



Augustinus Leyer 1735

(OCRopus: 97%, Tesseract: 82%)

TypeTable / Few Pages

Training Data

Type 1 109R (Palenomas_MUP)

Glyphe	Zeichen	Unicode
A	A	0041
B	B	0042
C	C	0043
D	D	0044
E	E	0045
F	F	0046
G	G	0047
H	H	0048
I	I	0049
K	K	004B
L	L	004C
M	M	004D

De inutilibus libris.
Inter precipuos pars est mihi reddita stultos
Prima:rego docili vastaq; vela manu.
En ego possideo multos/quos raro libellos
Perlego:tum lectos negligo:nec fapio:

De inutilibus libris.
Inter precipuos pars est mihi reddita stultos
Prima:rego docili vastaq; vela manu.
En ego possideo multos/quos raro libellos
Perlego:tum lectos negligo:nec fapio:

Tools: Aletheia and Franken+

Training Data

Roma fenes habuit consultoresq; peritos

Roma fenes habuit consultoresq; peritos

Nil sibi diuitię profunt:nil fortis auarę

Nil sibi diuitię profunt:nil fortis auarę

Optauit nemo sapiens/ in orbe potenter

Optauit nemo sapiens/ in orbe potenter

~50k to 100k



anyOCR - Automatic OCR Model

*anyOCR: A Combination of **OCRopus** and **tesseRECT** (**OCRoRECT**)*

- **Training Data:** only TypeTable
- **Higher Performance:**
 - Comparable to OCRopus
 - Independent to Noise & Character Seg.
 - Implicit Language Model



anyOCR - Automatic OCR Model

anyOCR: A Combination of *OCRopus* and *tesseRECT* (*OCRoRECT*)

- **Training Data:** only TypeTable
- **Higher Performance:**
 - Comparable to OCRopus
 - Independent to Noise & Character Seg.
 - Implicit Language Model

Type 1 109R (Palemonas_MUFI)

Glyphe	Zeichen	Unicode
A	A	0041
B	B	0042
C	C	0043
D	D	0044
E	E	0045
F	F	0046
G	G	0047
H	H	0048
I	I	0049
K	K	004B
L	L	004C
M	M	004D



anyOCR - Automatic OCR Model

anyOCR: A Combination of *OCROPUS* and *tesseRECT* (OCRoRECT)

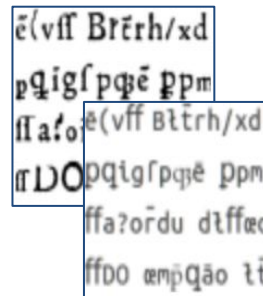
- **Training Data:** only TypeTable
- **Higher Performance:**
 - Comparable to OCROPUS
 - Independent to Noise & Character Seg.
 - Implicit Language Model

Type 1 109R (Palemonas_MUFI)

Glyphe	Zeichen	Unicode
A	A	0041
B	B	0042
C	C	0043
D	D	0044
E	E	0045
F	F	0046
G	G	0047
H	H	0048
I	I	0049
K	K	004B
L	L	004C
M	M	004D



Meaningless Text
Generator

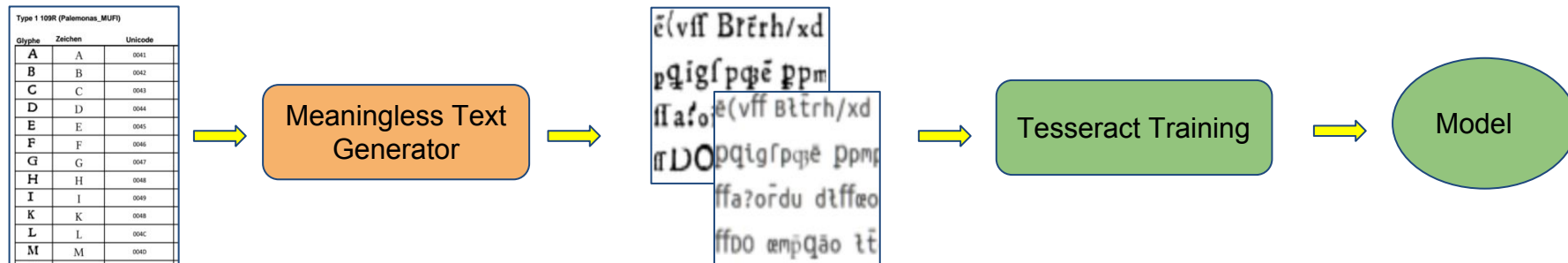




anyOCR - Automatic OCR Model

*anyOCR: A Combination of **OCROPUS** and **tesseRECT** (OCRoRECT)*

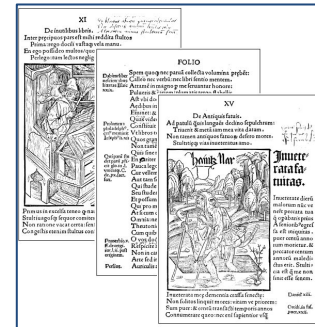
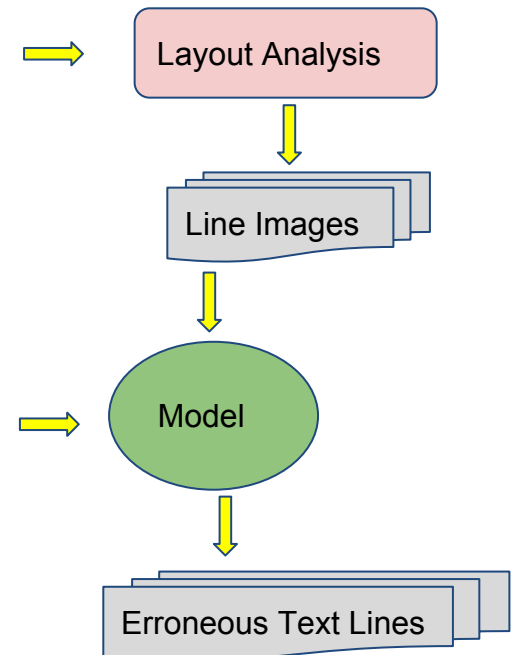
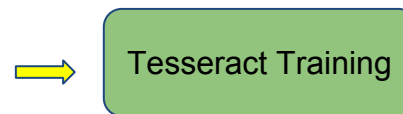
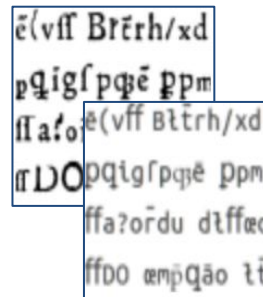
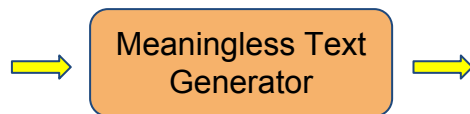
- **Training Data:** only TypeTable
- **Higher Performance:**
 - Comparable to OCROPUS
 - Independent to Noise & Character Seg.
 - Implicit Language Model



anyOCR: A Combination of OCRopus and tesseRECT (OCRoRECT)

- **Training Data:** only TypeTable
- **Higher Performance:**
 - Comparable to OCRopus
 - Independent to Noise & Character Seg.
 - Implicit Language Model

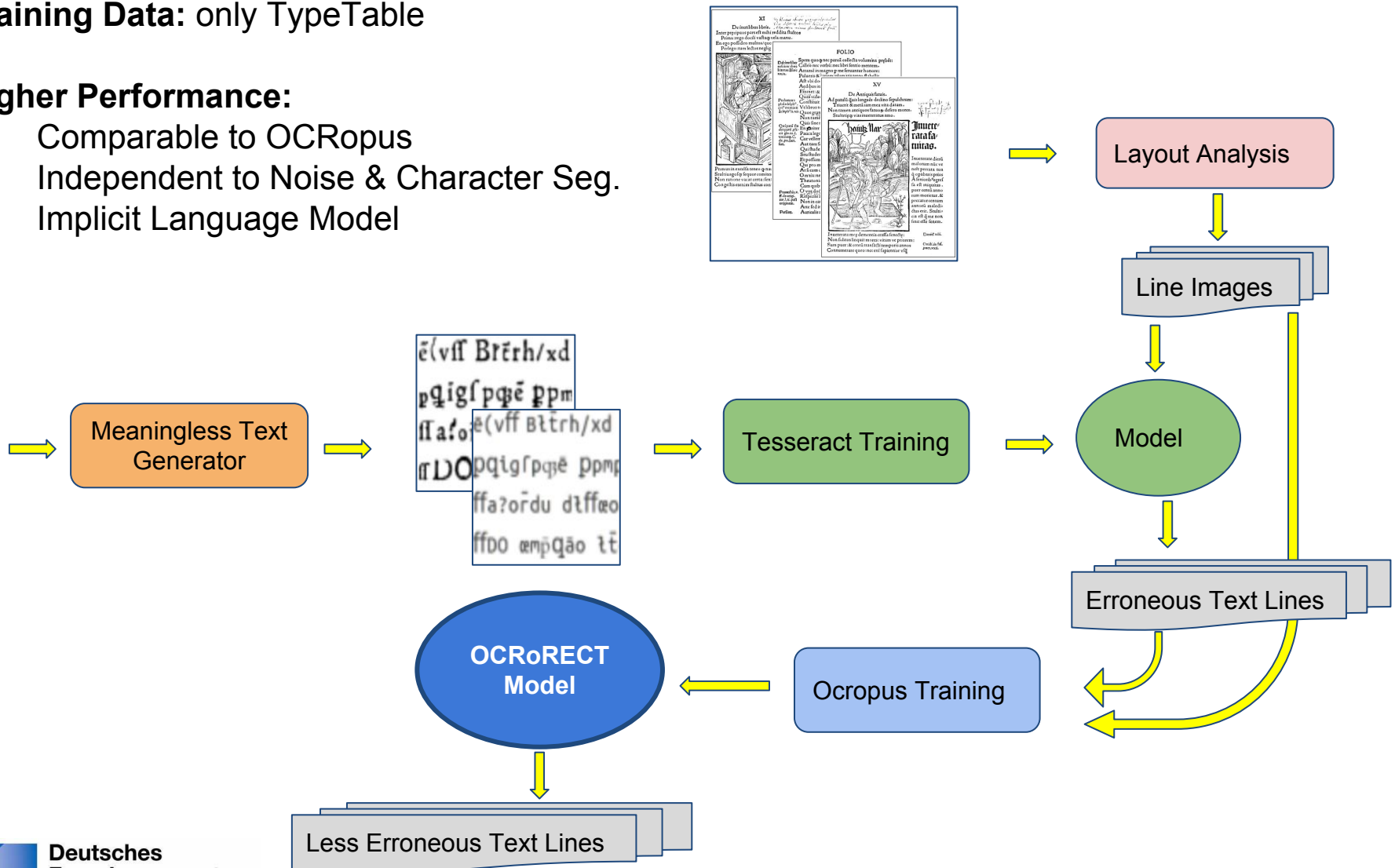
Type 1 109R (Palemonas_MUFI)		
Glyphe	Zeichen	Unicode
A	A	0041
B	B	0042
C	C	0043
D	D	0044
E	E	0045
F	F	0046
G	G	0047
H	H	0048
I	I	0049
K	K	004B
L	L	004C
M	M	004D



anyOCR: A Combination of OCRopus and tesseRECT (OCRoRECT)

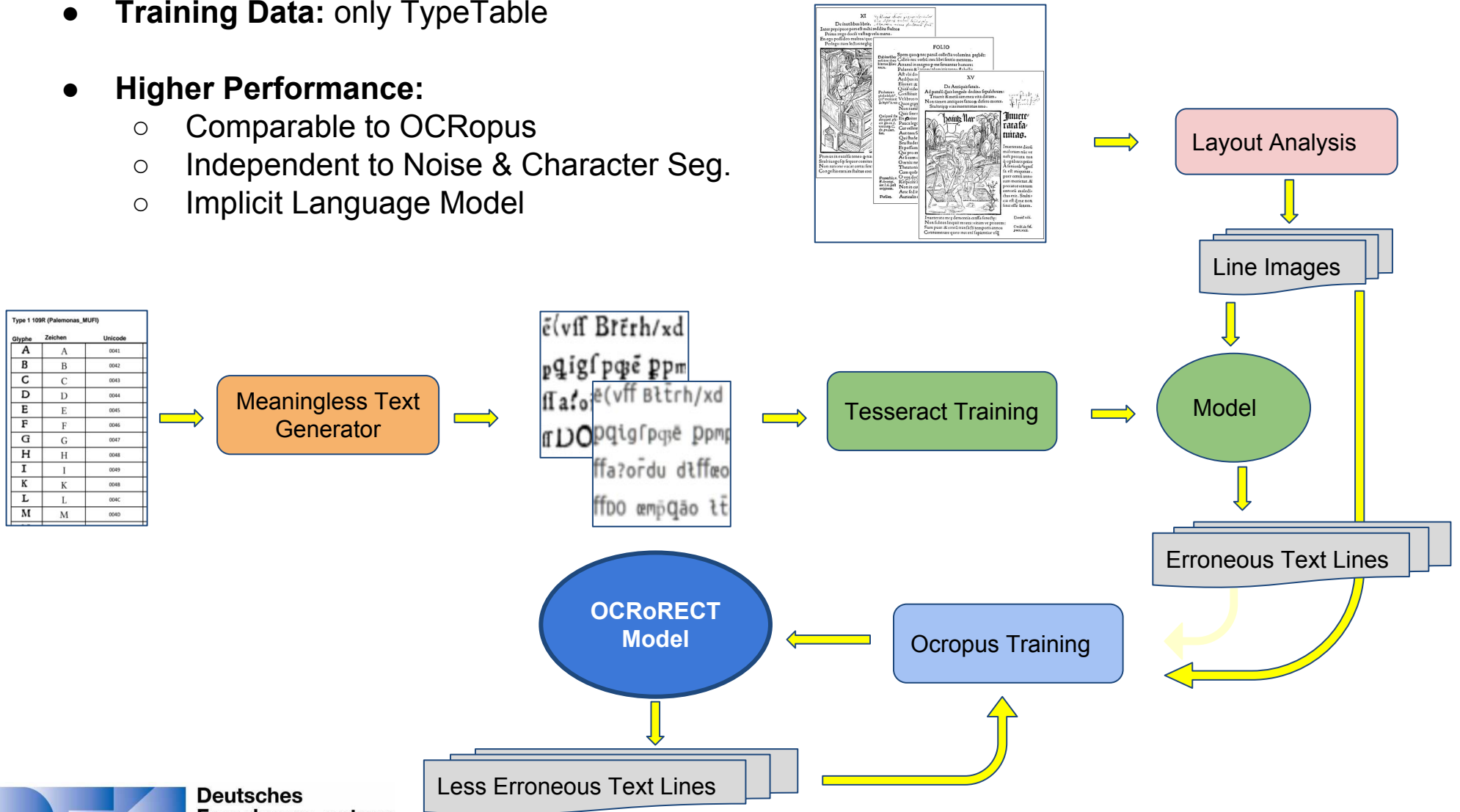
- **Training Data:** only TypeTable
- **Higher Performance:**
 - Comparable to OCRopus
 - Independent to Noise & Character Seg.
 - Implicit Language Model

Type 1109R (Palemonas_MUFI)		
Glyphe	Zeichen	Unicode
A	A	0041
B	B	0042
C	C	0043
D	D	0044
E	E	0045
F	F	0046
G	G	0047
H	H	0048
I	I	0049
K	K	004B
L	L	004C
M	M	004D



anyOCR: A Combination of OCRopus and tesseRECT (OCRoRECT)

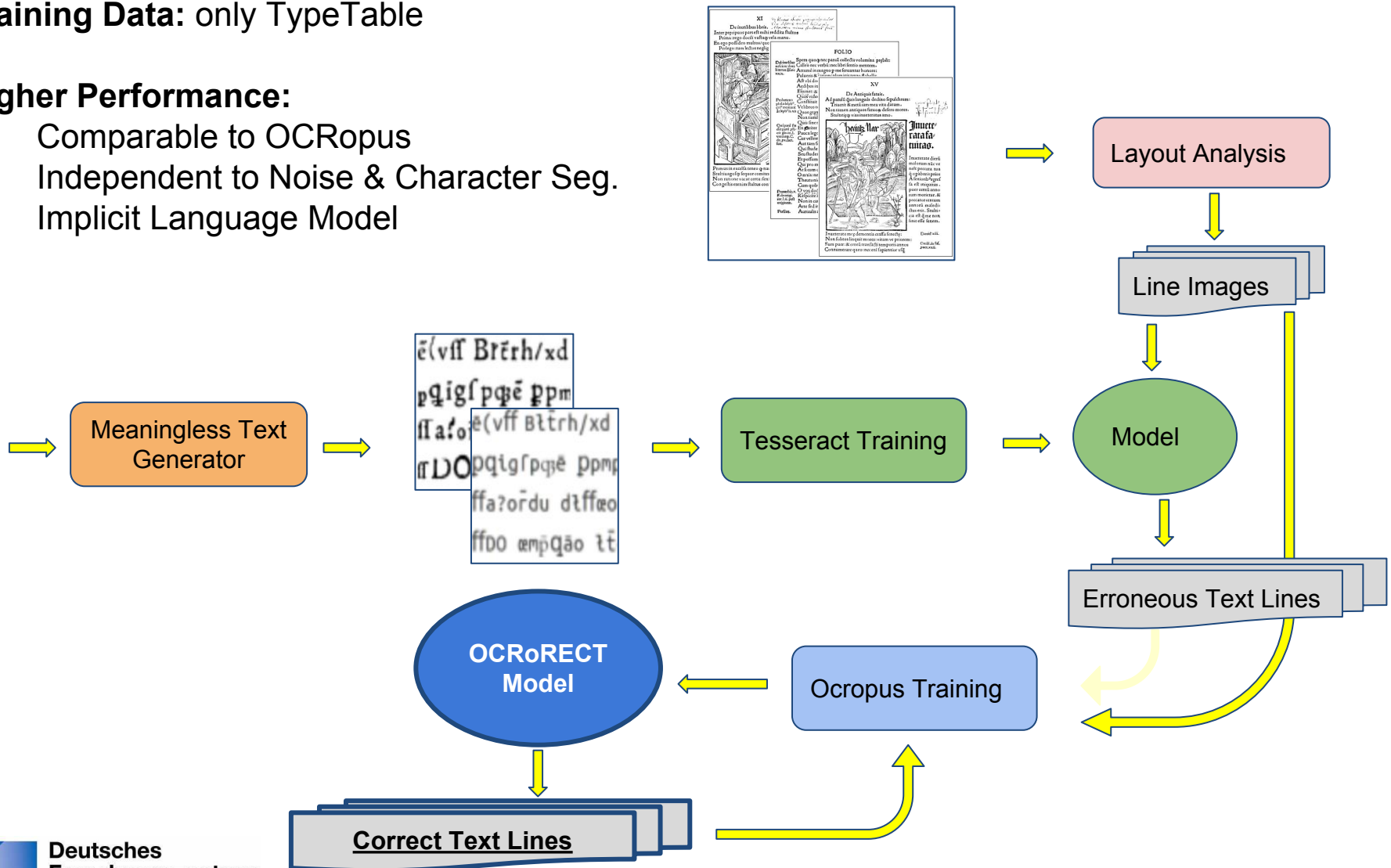
- **Training Data:** only TypeTable
- **Higher Performance:**
 - Comparable to OCRopus
 - Independent to Noise & Character Seg.
 - Implicit Language Model



anyOCR: A Combination of OCROPUS and tesseRECT (OCRoRECT)

- **Training Data:** only TypeTable
- **Higher Performance:**
 - Comparable to OCROPUS
 - Independent to Noise & Character Seg.
 - Implicit Language Model

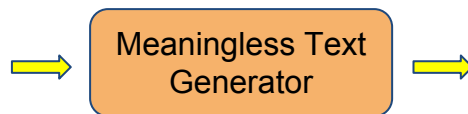
Type 1 100R (Palemonas_MUFI)		
Glyphe	Zeichen	Unicode
A	A	0041
B	B	0042
C	C	0043
D	D	0044
E	E	0045
F	F	0046
G	G	0047
H	H	0048
I	I	0049
K	K	004B
L	L	004C
M	M	004D



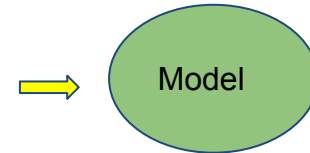
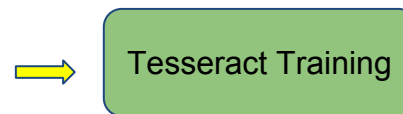
anyOCR: A Combination of OCRopus and tesseRECT (OCRoRECT)

- **Training Data:** only TypeTable
- **Higher Performance:**
 - Comparable to OCRopus
 - Independent to Noise & Character Seg.
 - Implicit Language Model

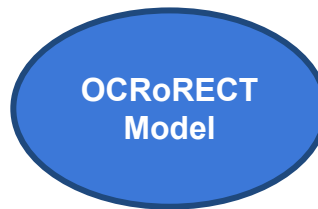
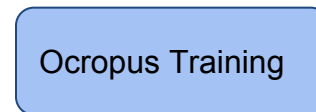
Type 1 100R (Palemonas_MUFI)		
Glyphe	Zeichen	Unicode
A	A	0041
B	B	0042
C	C	0043
D	D	0044
E	E	0045
F	F	0046
G	G	0047
H	H	0048
I	I	0049
K	K	004B
L	L	004C
M	M	004D



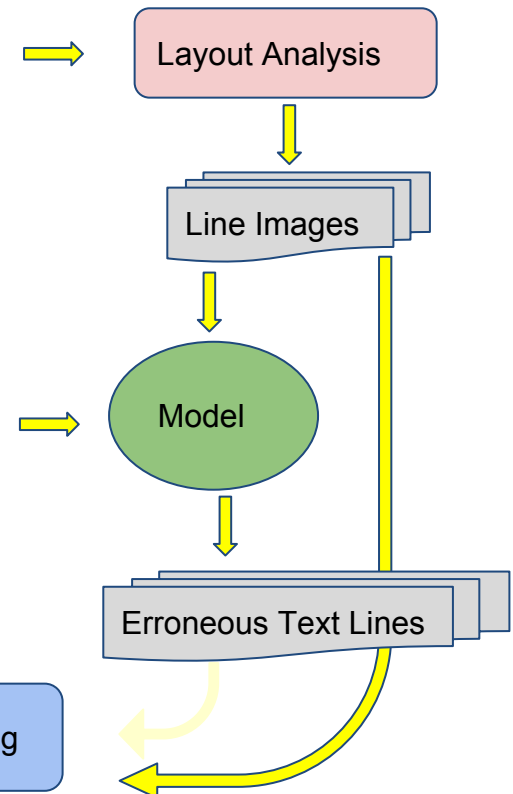
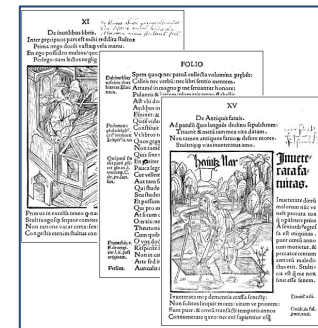
ē(vff Btrrh/xd
p q i g s p q ē p p m
ff a? o ē(vff Btrrh/xd
ff a? o r d u d i f f e o
ff d o æ p q ā o i t



Erroneous Text Lines



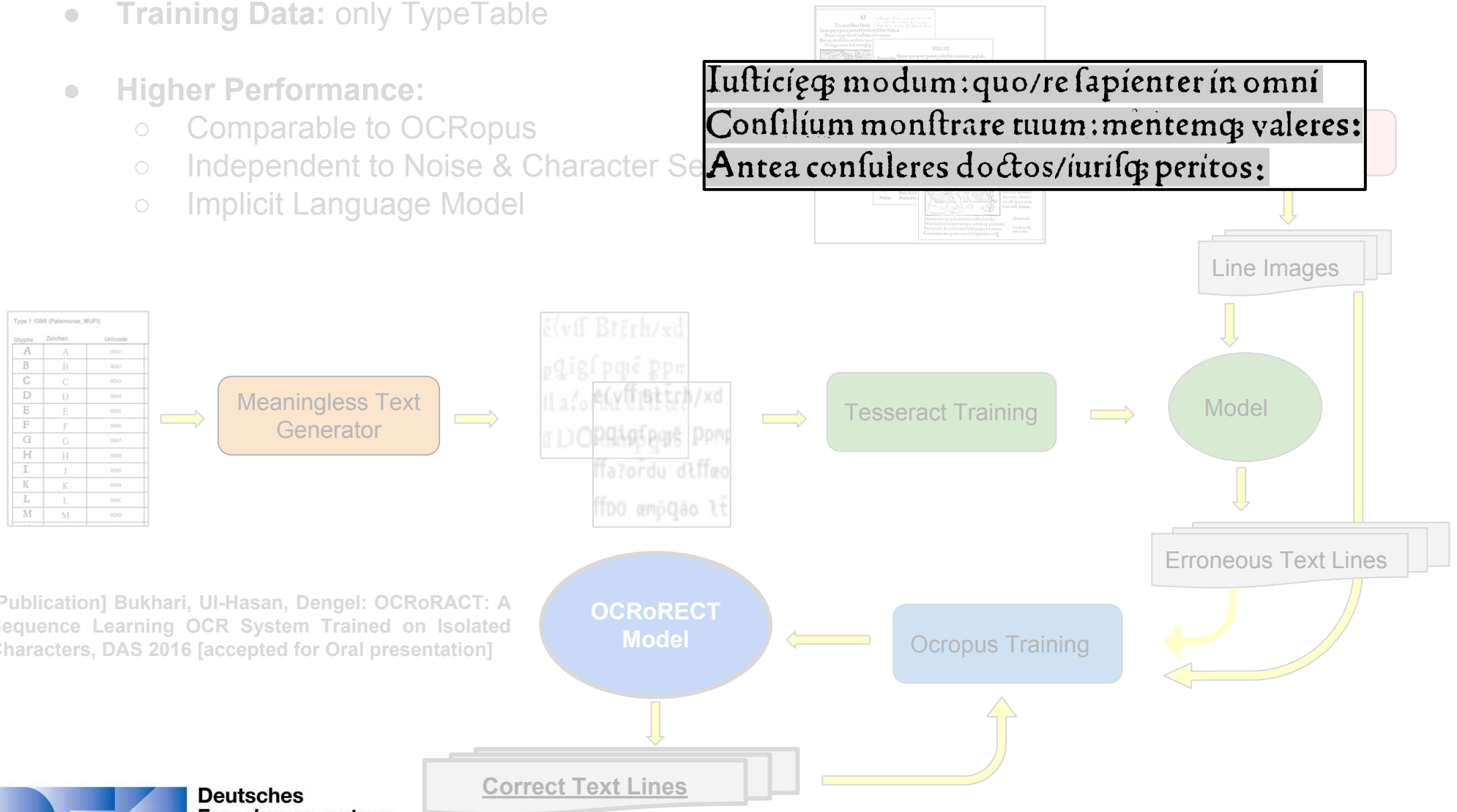
Correct Text Lines



[Publication] Bukhari, UI-Hasan, Dengel: OCRoRACT: A Sequence Learning OCR System Trained on Isolated Characters, DAS 2016 [accepted for Oral presentation]

anyOCR: A Combination of OCROPUS and tesseRECT (OCRoRECT)

- Training Data: only TypeTable
- Higher Performance:
 - Comparable to OCROPUS
 - Independent to Noise & Character Set
 - Implicit Language Model



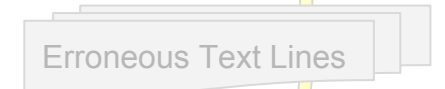
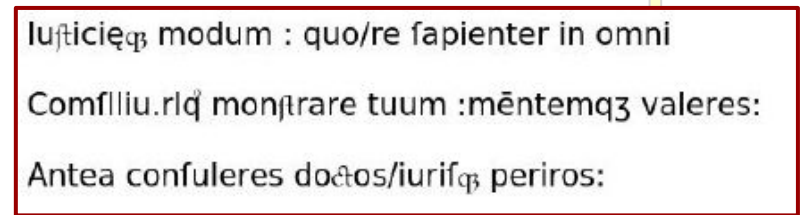
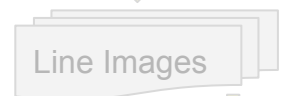
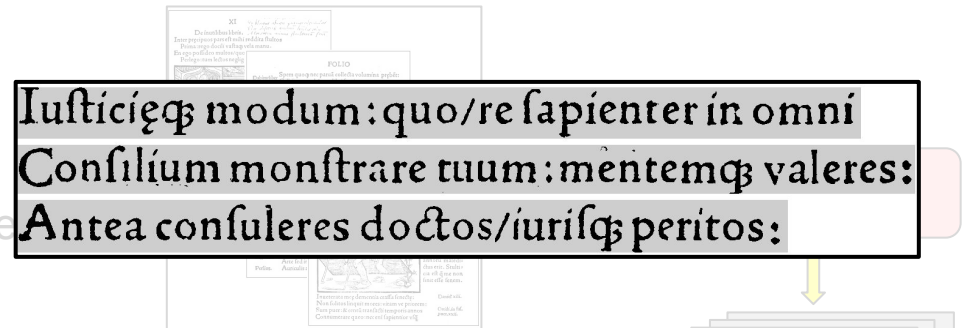
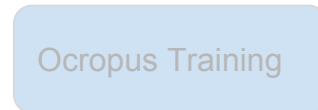
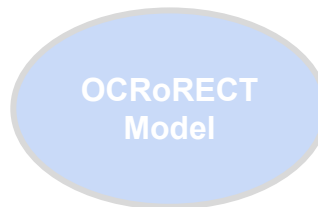
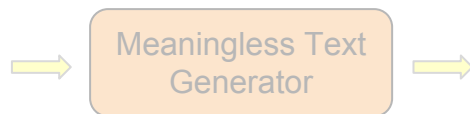
[Publication] Bukhari, UI-Hasan, Dengel: OCRoRACT: A Sequence Learning OCR System Trained on Isolated Characters, DAS 2016 [accepted for Oral presentation]

anyOCR: A Combination of OCROPUS and tesseRECT (OCRoRECT)

- Training Data: only TypeTable
- Higher Performance:
 - Comparable to OCROPUS
 - Independent to Noise & Character Set
 - Implicit Language Model

Type 1109R (Palemonas_MUFI)

Glyphe	Zeichen	Unicode
A	A	0041
B	B	0042
C	C	0043
D	D	0044
E	E	0045
F	F	0046
G	G	0047
H	H	0048
I	I	0049
K	K	004A
L	L	004B
M	M	004C



[Publication] Bukhari, UI-Hasan, Dengel: OCRoRACT: A Sequence Learning OCR System Trained on Isolated Characters, DAS 2016 [accepted for Oral presentation]

anyOCR: A Combination of OCROPUS and tesseRECT (OCRoRECT)

- **Training Data:** only TypeTable
- **Higher Performance:**
 - Comparable to OCROPUS
 - Independent to Noise & Character Sequence
 - Implicit Language Model

Type 1 100R (Palemonas_MUFI)		
Glyphe	Zeichen	Unicode
A	A	0041
B	B	0042
C	C	0043
D	D	0044
E	E	0045
F	F	0046
G	G	0047
H	H	0048
I	I	0049
K	K	004A
L	L	004B
M	M	004C

Meaningless Text Generator



Iusticięqꝫ modum:quo/re sapienter in omni
Conflum monſtrare tuum: mentemqꝫ valeres:
Antea confuleres doctos/iurifqꝫ peritos:

RECT

Ocropus Training

Correct Text Lines

Iusticięqꝫ modum:quo/re sapienter in omni
Conſilium monſtrare tuum: mētemqꝫ valeres:
Antea confuleres doctos/iurifqꝫ peritos:

Line Images

Iuſticięqꝫ modum : quo/re ſapienter in omni
Comſiliu.rlq monſtrare tuum :mētemqꝫ valeres:
Antea confuleres doctos/iurifqꝫ periros:

Erroneous Text Lines

[Publication]
Sequence L
Characters,

Qualitative Performance Evaluation of **OCRopus++**

Iusticięqꝫ modum: quo/re sapienter in omni
Consilium monstrare tuum: mentemqꝫ valeres:
Antea consuleres doctos/iurifqꝫ peritos:

Iuſticięqꝫ modum : quo/re ſapienter in omni
Comſilliu.rlqꝫ monſtrare tuum :mētemqꝫ valeres:
Antea conſuleres doctos/iurifqꝫ periros:

Tesseract

Iuſticięqꝫ modum:quo/re ſapienter in omni
Conſilium monſtrare tuum: mentemqꝫ valeres:
Antea conſuleres doctos/iurifqꝫ peritos:

anyOCR: OCRoRECT

Qualitative Performance Evaluation of **OCRopus++**

Iusticięqꝫ modum:quo/re sapienter in omni
Consilium monſtrare tuum:mentemqꝫ valeres:
Antea conſuleres doctos/iurifqꝫ peritos:

Iuſticięqꝫ modum : quo/re ſapienter in omni
Comſilli.riqꝫ monſtrare tuum :mētemqꝫ valeres:
Antea conſuleres doctos/iurifqꝫ periros:

Tesseract

Iuſticięqꝫ modum:quo/re ſapienter in omni
Conſilium monſtrare tuum:mentemqꝫ valeres:
Antea conſuleres doctos/iurifqꝫ peritos:

OCRopus

Iuſticięqꝫ modum:quo/re ſapienter in omni
Conſilium monſtrare tuum: mentemqꝫ valeres:
Antea conſuleres doctos/iurifqꝫ peritos:

anyOCR: OCRoRECT

Performance Evaluation of **OCROPUS++**:

OCR Model Dataset	Tesseract	OCROPUS	anyOCR - OCRoRECT
Adam von Bodenstein 1557	78% [1]	99% [1]	-
Augustinus Leyer 1735	82% [1]	97% [1]	-

[1] Springmann, Ocrocis Tutorial: cistern.cis.lmu.de/ocrocis/tutorial.pdf

Performance Evaluation of **OCROPUS++**:

OCR Model Dataset	Tesseract	OCROPUS	anyOCR - OCRoRECT
Adam von Bodenstein 1557	78% [1]	99% [1]	-
Augustinus Leyer 1735	82% [1]	97% [1]	-
Basel 1497 (Narrenschif, Latin)	77% [2]	98% [2]	

[1] Springmann, Ocrocis Tutorial: cistern.cis.lmu.de/ocrocis/tutorial.pdf

[2] Bukhari, UI-Hasan, Dengel: OCRoRACT: A Sequence Learning OCR System Trained on Isolated Characters, DAS 2016 [accepted for Oral presentation]

Performance Evaluation of **OCROPUS++**:

OCR Model Dataset	Tesseract	OCROPUS	anyOCR - OCRoRECT
Adam von Bodenstein 1557	78% [1]	99% [1]	-
Augustinus Leyer 1735	82% [1]	97% [1]	-
Basel 1497 (Narrenschif, Latin)	77% [2]	98% [2]	95% [2]

[1] Springmann, Ocrocis Tutorial: cistern.cis.lmu.de/ocrocis/tutorial.pdf

[2] Bukhari, UI-Hasan, Dengel: OCRoRACT: A Sequence Learning OCR System Trained on Isolated Characters, DAS 2016 [accepted for Oral presentation]

Conclusion

Conclusion

- Layout Analysis is an open-challenging problem for complex documents
- OCR recognizers those need less training data fail to achieve good performance (e.g. Tesseract) and those produce better results require a lot of training data (e.g. OCRopus)

Conclusion

- Layout Analysis is an open-challenging problem for complex documents
- OCR recognizers those need less training data fail to achieve good performance (e.g. Tesseract) and those produce better results require a lot of training data (e.g. OCRopus)
- Presented **OCROPUS++** OCR System
 - advanced layout analysis
 - anyOCR - OCRoRECT recognizer

Conclusion

- Layout Analysis is an open-challenging problem for complex documents
- OCR recognizers those need less training data fail to achieve good performance (e.g. Tesseract) and those produce better results require a lot of training data (e.g. OCRopus)
- Presented **OCROPUS++** OCR System
 - advanced layout analysis
 - anyOCR - OCRoRECT recognizer
- In future, add more features in **OCROPUS++**

Conclusion

- Layout Analysis is an open-challenging problem for complex documents
- OCR recognizers those need less training data fail to achieve good performance (e.g. Tesseract) and those produce better results require a lot of training data (e.g. OCRopus)
- Presented **OCROPUS++** OCR System
 - advanced layout analysis
 - anyOCR - OCRoRECT recognizer
- In future, add more features in **OCROPUS++**
- Publications:
 - Bukhari, UI-Hasan, Dengel, "OCRoRACT: A Sequence Learning OCR System Trained on Isolated Characters", DAS 2016
 - Bukhari, UI-Hasan, Dengel, "Meaningless Text OCR Model for Medieval Scripts", 2nd International Conference on Natural Sciences and Technology in Manuscript Analysis 2016, Germany.
 - Bukhari, Jenckel, Dengel, "Clustering Benchmark for Characters in Historical Documents", DAS 2016.
 - Bukhari, Nunamaker, Borth, Dengel, "A Tesseract Based OCR Framework For Historical Document Lacking Ground-Truth", ICIP 2016 [Under Review]

Thank you, ... Questions?



Address:

Dr.-Ing. Syed Saqib Bukhari

DFKI GmbH

P.O. Box 2080

D-67608 Kaiserslautern

email: saqib.bukhari@dfki.de

<http://www.dfki.de/~bukhari>

